# FLAVR: Flow-Agnostic Video Representations for Fast Frame Interpolation

Tarun Kalluri *
UCSD

Deepak Pathak
CMU

Manmohan Chandraker
UCSD

Du Tran
Facebook AI

https://tarun005.github.io/FLAVR/

## Abstract

*A majority of methods for video frame interpolation compute bidirectional optical flow between adjacent frames of a video, followed by a suitable warping algorithm to generate the output frames. However, approaches relying on optical flow often fail to model occlusions and complex non-linear motions directly from the video and introduce additional bottlenecks unsuitable for widespread deployment. We address these limitations with FLAVR, a flexible and efficient architecture that uses 3D space-time convolutions to enable end-to-end learning and inference for video frame interpolation. Our method efficiently learns to reason about non-linear motions, complex occlusions and temporal abstractions, resulting in improved performance on video interpolation, while requiring no additional inputs in the form of optical flow or depth maps. Due to its simplicity, FLAVR can deliver 3× faster inference speed compared to the current most accurate method on multi-frame interpolation without losing interpolation accuracy. In addition, we evaluate FLAVR on a wide range of challenging settings and consistently demonstrate superior qualitative and quantitative results compared with prior methods on various popular benchmarks including Vimeo-90K, UCF101, DAVIS, Adobe, and GoPro. Finally, we demonstrate that FLAVR for video frame interpolation can serve as a useful self-supervised pretext task for action recognition, optical flow estimation, and motion magnification.*

## 1. Introduction

Video frame interpolation [35, 1, 77, 27, 22, 43, 40, 32, 9, 41] is a challenging problem in video analysis that aims to overcome the limited acquisition frame rate and exposure time of commercial video cameras, where the task is to generate non-existent intermediate frames between existing ones. A large number of prior works use *bidirectional flow warping* for frame interpolation [22, 77, 41], where the input frames are used to estimate bidirectional optical
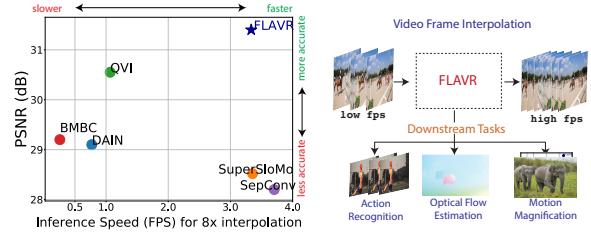
Figure 1: **Our contributions** We propose FLAVR, a simple and efficient architecture for single shot multi-frame interpolation. The plot of accuracy (PSNR) vs. inference speed (fps) of FLAVR compared with current methods on GoPro 8x interpolation with 512×512 input images. FLAVR offers optimum trade off between accuracy and inference speed. Also, FLAVR is a useful self-supervised pretext task for various downstream applications.

flow maps from a reliable flow estimator network, possibly along with additional information like monocular depth maps [1] and occlusion masks [2]. The interpolated frames at intermediate time steps are then generated either by using backward [22, 1] or forward warping [40, 41]. However, the optical flow-based approaches, as well as proposed alternatives [43, 42, 49, 6, 27], have to confront one or more of the following limitations: (1) **Computational Costs** As they rely on optical flow and pixel level warping procedures, they are inefficient at both training and inference making them less suitable for end applications. For example, [76], [1] and [47] take order of seconds to generate frames for 8×interpolation while also requiring users to deploy custom cuda kernels that prohibit seamless deployment across edge devices. (2) **Modeling Complex Trajectories** The modeling capacity is limited to account for only linear [22, 1] or quadratic [76, 7] motion trajectories, and extending these to account for more complex motions is non-trivial. (3) **Representation Inflexibility** By accepting pre-computed optical flows as inputs, they focus on learning only spatial warping and interpolation, thus the representations learned in the process are not useful beyond frame interpolation.

In this work, we aim to bring a smooth tradeoff between visual quality and inference speed for video interpolation. We do so by proposing FLAVR (Flow-Agnostic Video Representation network), which jointly addresses the aforemen-

tioned limitations. FLAVR is a simple yet efficient CNN architecture that utilizes spatio-temporal convolutions for predicting intermediate frames of a video. Without demanding access to external flow or depth maps, our model is able to make single-shot, end-to-end multiple-frame predictions. It naturally handles complex motions and occlusions through learning from large scale video data, significantly improves inference speed and ease of deployment compared to prior approaches, while achieving state-of-the art interpolation accuracy (Figure. 1).

We also posit that models learned from video should be able to simultaneously reason about intricate synergy between objects, motions and actions for accurate frame interpolation. This is because different actions and objects have different motion signatures, and it is essential to precisely capture these properties through the representations learned for accurate frame interpolation. We ground this argument in the context of self-supervised representation learning [10, 48, 65, 16] from videos. While popular pretext tasks for unsupervised learning from video include predicting ordering of frames [73, 28, 12, 38, 67], tracking color [63, 66], correspondence across time [21, 20] or contrastive predictive coding [14, 15], we take a fundamentally different approach and use frame interpolation as a pretext task to learn spatio-temporal representations from large scale unlabeled video. We demonstrate the validity of the approach by improving performance on action recognition and optical flow estimation tasks compared with the training from scratch baseline as well as other self-supervised approaches based on pixel level pretext tasks.

In summary, we make the following contributions:

- We propose FLAVR, a scalable, flow-free, efficient 3D CNN architecture for video frame interpolation. To the best of our knowledge, FLAVR is the first video frame interpolation approach that is **both** *optical flow-free* and able to make *single-shot multiple-frame predictions* (Section. 3).
- FLAVR is quantitatively and qualitatively superior or comparable to current approaches on multiple standard benchmarks including Vimeo-90K, UCF101, DAVIS, Adobe and GoPro while offering the best trade-off in terms of accuracy and inference speed for video interpolation (Figure. 1 and Figure. 3a)[1].
- We demonstrate that video representations self-supervisedly learned by FLAVR can be reused for various downstream tasks such as action recognition, optical flow estimation, and motion magnification (Section. 6).

---

[1]A previous version of FLAVR erroneously reported faster runtime for $8\times$ interpolation. We corrected the bugs and reported the accurate runtime here.

## 2. Related Work

**Video Frame Interpolation** Video frame interpolation is a classical computer vision problem [34] and recent methods take one of phase based, flow based, or kernel based approaches.

*Phase based models* for video interpolation [36, 35] consider each frame as a linear combination of wavelets. The phase and magnitude of each phase is then interpolated using classical interpolation [36] or deep learning based algorithms [36] across multi-scale pyramid levels.

*Optical Flow based methods* [22, 1, 76, 77, 7, 31, 2, 79, 40, 78, 80, 41] use a optical flow prediction network, *e.g.* PWC-Net [59], to compute bidirectional optical flow between the input frames [22] that guides frame synthesis along with occlusion masks [22, 2, 77] or monocular depth maps [1] to reason about occlusions. Contextual warping [40], softmax splatting [41], cycle constraints [51, 31] and meta learning [8] have also shown to be effective tools in improving the performance of such methods. While being largely successful in generating realistic intermediate frames, their performance is limited by the accuracy of the underlying flow estimator, which can be noisy in presence of large motions and complex occlusions resulting in outputs that have noticeable artifacts. They also assume uniform linear motion between the frames which is far from ideal for real world videos. Recent works relax the assumption to propose quadratic warping by taking a larger input context [76, 7, 30] at the cost of increased model complexity or inference time.

*Kernel based methods*, on the other hand, skip flow computation altogether by predicting spatially adaptive filters to resample from input frames [43, 42, 49, 6, 32, 55]. Most works only consider a local neighborhood patch for resampling, and using larger kernels to increase receptive field incurs significant memory and inference time overhead. CAIN [9] uses channel attention as suitable ingredients for frame interpolation but fails to capture complex spatio temporal dependencies explicitly between input frames. We address all these issues in this work by designing an end to end architecture that directly predicts intermediate frames from a given video by learning to reason motion trajectories and properties through 3D space time convolutions while jointly optimizing for output quality and inference time.

A key distinction with works like temporal super resolution (TSR) [81, 13, 54, 71, 56] and deep video deblurring (DVD) [58, 69] is that that TSR and DVD are focused on simultaneously improving the spatial *and* temporal resolution. In this work, we only consider temporal frame interpolation and leave possible extensions of FLAVR to TSR or DVD as a future work.

**Spatio-temporal Filtering** Due to their proven success in capturing complex spatial and temporal dependencies, 3D space-time convolutions are very commonly used in video understanding tasks like action recognition [60, 61, 72, 4,
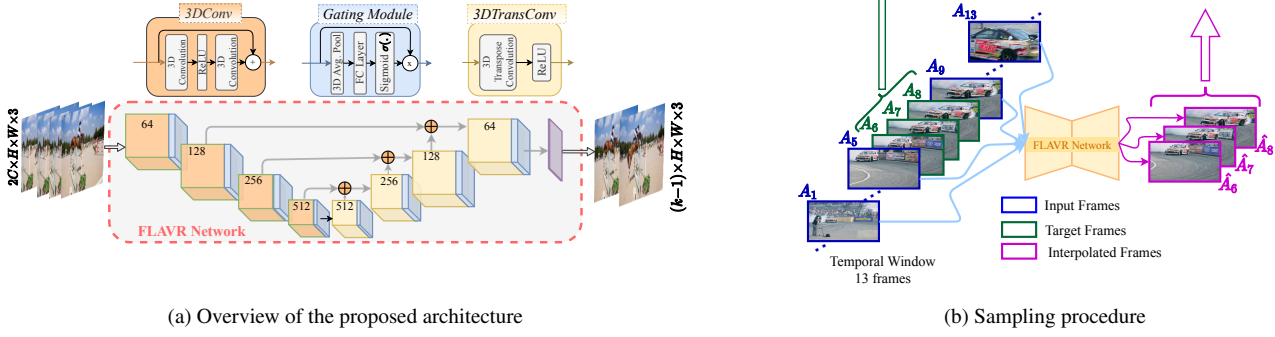
(a) Overview of the proposed architecture

(b) Sampling procedure

Figure 2: **FLAVR Architecture**. (a) Our FLAVR is U-Net style architecture with 3D space-time convolutions (orange blocks) and deconvolutions (yellow blocks). We use channel gating after all (de-)convolution layers (blue blocks). The final prediction layer (the purple block) is implemented as a convolution layer to project the 3D feature maps into $(k-1)$ frame predictions. This design allows FLAVR to predict multiple frames in one inference forward pass. (b) A concrete example of our sampling procedure for $4\times$ interpolation ($k=4$) with 4-frame input ($C=2$). Best viewed in color.

11], action detection [57, 74], and captioning [75]. We explore the use of 3D convolutions for the problem of temporal frame interpolation which requires modeling complex temporal abstractions between inputs for generating accurate and sharp predictions.

## 3. Frame Interpolation using FLAVR

The problem of video frame interpolation is to generate a high frame-rate video from a lower frame-rate input video. With $k$ as the *interpolation factor*, the $k\times$-video frame interpolation problem is to predict $(k-1)$ additional intermediate frames for every original frame of the input video that are both spatially and temporally consistent with the rest of the video. We now present our approach to solve the task, and explain the details of the proposed architecture.

**Sampling Training Data from Unlabeled Videos** We can directly generate inputs and ground truths required for training from raw videos as follows. Let $k$ be the interpolation factor, $V$ is the original video with a frame rate $f$ FPS. In order to generate training data for the $k\times$-video frame interpolation problem, we sub-sample frames of $V$ with a sampling stride of $k$ to form a low frame rate video $\bar{V}$ with $\frac{f}{k}$ fps. Then, to perform interpolation between any two frames $A_i$, $A_{i+1}$ of $\bar{V}$, we use a temporal window of size $2C$ centered around $A_i$ and $A_{i+1}$ as the input, and all frames between $A_i$ and $A_{i+1}$ in original video $V$ as the ground truth. Our network is flexible to handle any temporal context instead of just $A_i$, $A_{i+1}$ [22, 1], which helps us to model complex trajectories and improve interpolation accuracy. The sampled input frames are concatenated in the temporal dimension resulting in input dimension $2C\times H\times W\times 3$, where $H, W$ are the spatial dimensions of the input video.

An illustration of this sampling procedure is demonstrated in Figure. 2b for the case of $4\times$ interpolation ($k=4$) with two context inputs from the past and future ($C = 2$). In

this case, the frames $\{A_1, A_5, A_9, A_{13}\}$ are used as inputs to predict the 3 intermediate frames of $\{A_6, A_7, A_8\}$. Intuitively, the frames in the immediate neighborhood would be more relevant for frame interpolation than frames farther out. In our experiments, we find that for most common settings, using four context frames ($C = 2$) is sufficient for accurate prediction. We present a detailed analysis of the effect of input context in Section. 5.2.

**Architecture Overview** We present the proposed architecture of FLAVR in Figure. 2a. FLAVR is a 3D U-Net obtained by extending the popular 2D Unet [52] used in pixel generation tasks, by replacing all the 2D convolutions in the encoder and decoder with 3D convolutions (*3DConv*) to accurately model the temporal dynamics between the input frames, invariably resulting in better interpolation quality. Each 3D filter is a 5-dimensional filter of size $c_i\times c_o\times t\times h\times w$, where $t$ is the temporal size and $(h, w)$ is the spatial size of the kernel. $c_i$ and $c_o$ are the number of input and output channels in the layer. The additional temporal dimension is useful in modeling the temporal abstractions like motion trajectories, actions or correspondences between frames in the video. We show in Section. 6 that our network indeed learns useful representations along the temporal dimensions that can be reused in downstream tasks like action recognition with limited labeled data.

Practically any 3D CNN architecture can be used as the encoder, and we use ResNet-3D (R3D) with 18 layers [61] in our work to strike a balance between complexity and accuracy. We remove the last classification layer from R3D-18, resulting in 5 conv blocks *conv1* to *conv5*, each made up of two 3D convolutional layers and a skip connection. We also remove all temporal striding, as downsampling operations like striding and pooling are known to remove details that are crucial for generating sharper images. However, we do use spatial stride of 2 in *conv1*, *conv3* and *conv4* blocks of the network to keep the computation manageable.

The decoder essentially constructs the output frames from a deep latent representation by progressive, multi-scale feature upsampling and feature fusion from the encoder. For upsampling, we use 3D transpose convolution layers (*3DTransConv*) with a stride of 2. To handle the commonly observed checkerboard artefacts [45], we add a *3DConv* layer after the last *3DTransConv* layer. We also include skip connections that directly combine encoder features with the corresponding decoder along the channels to fuse the low level and high level information necessary for accurate and sharp interpolation.

The output of the decoder, which is a 3D feature map, is then passed through a temporal fusion layer, implemented by a 2D conv, in which the features from the temporal dimension are concatenated along the channels resulting in a 2D spatial feature map. This helps to aggregate and merge information present in multiple frames for prediction. Finally, this output is passed through a $7\times7$ 2D convolution kernel that predicts output of size $H\times W\times3(k-1)$, which is then split along the channel dimension to get the $(k-1)$ output frames.

**Spatio-Temporal Feature Gating** Feature gating technique is used as a form of self-attention mechanism in deep neural networks for action recognition [37, 72], image classification [18] and video interpolation [9]. We apply the gating module after every layer in our architecture. Given an intermediate feature dimension of size $f_i = C\times T\times H\times W$, the output $f_o$ of the gating layer is given by $f_o = \sigma(W.pool(f_i) + b) \odot f_i$ where $W \in \mathbb{R}^{C\times C}$ and $b \in \mathbb{R}^C$ are learnable weight and bias parameters and *pool* is a spatio-temporal pooling layer. Such a feature gating mechanism would suitably learn to upweight certain relevant dimensions of the feature maps that learn useful cues for frame interpolation, like motion boundaries (see Section. 5.2).

**Loss Function** We train the whole network end to end using the following L1 loss:

$$\mathcal{L}(\{\hat{I}\}, \{I\}) = \frac{1}{N} \sum_{i=1}^{N} \sum_{j=1}^{k-1} ||\hat{I}_j^{(i)} - I_j^{(i)}||_1 \qquad (1)$$

where $\{\hat{I}_j^{(i)}\}$ and $\{I_j^{(i)}\}$ are the j-th predicted and the j-th ground truth frame of the $i^{\text{th}}$ target sample, $k$ is the interpolation factor, and $N$ is the size of the mini-batch used in training.

**Representation Learning using FLAVR** In order to successfully predict intermediate frames, it is essential for FLAVR to accurately reason about motion trajectories, estimate and capture motion patterns specific to objects, and reconstruct both high level semantic detail and low level texture details. It is interesting to understand what types of motion information the networks learned and which tasks this representation is useful for. Therefore, we examine the possibility of using video frame interpolation in the context of unsupervised representation learning by pre-training FLAVR on the task of frame interpolation, and reusing the learned feature representations for the tasks of action recognition, optical flow estimation, and motion magnification.

## 4. Experimental Setup

**Datasets**. We use Vimeo-90K dataset [77] which consists of $91,701$ frame septuplets extracted from 30FPS videos for training single frame interpolation networks. We train our model on the train split and evaluate it on the test split of the dataset. Following [76], we additionally verify the *generalization* capability of our proposed approach. For single frame interpolation, we report the performance on the 100 quintuples generated from UCF101[24] and 2,847 quintuples generated from DAVIS dataset [50]. For multi frame interpolation, we use GoPro [39] as the training set, and report results on the Adobe dataset [58] and GoPro dataset[39] for $4\times$ and $8\times$ interpolations. For fair comparison, we follow the same settings used in [76] to report the result on all these datasets.

**Training Details**. We use the sampling strategy discussed in Section. 3 for generating inputs during training and testing. For data augmentation, we exploit the symmetry of the problem by randomly selecting input sequences during training and inverting the temporal order of the frames. In addition, we also horizontally flip all frames of randomly selected inputs. We train the network for 200 epochs using Adam optimizer [25] with an initial learning rate of $2\times10^{-4}$. We use a mini-batch size of $64$ on Vimeo-90K dataset and $32$ on GoPro dataset, and train our network on 8 2080Ti GPUs. We reduce the learning rate by half whenever the training plateaus which is cross-validated by the validation set. We apply mean normalization once for every mini-batch of input frames separately rather than using global mean normalization or batch normalization inside the network to achieve training stability. We provide more details about the training procedure in Appendix G.

**Evaluation Metrics**. Following previous works, we use PSNR and SSIM metrics to report the quantitative results of our method. For multi-frame interpolation we report the average value of the metric over all the predicted frames. Since these quantitative measures do not strongly correlate with the human visual system [44], we also conduct a user study to analyze and compare our generated videos with other competing approaches.

**Baselines**. We perform comparisons with the following baselines that perform single and multi frame video interpolation. (i) **DAIN [1]** performs depth aware frame interpolation using the depth and flow maps as inputs for occlusion reasoning in addition to the RGB frames. (ii) **QVI [76]** computes quadratic flow prediction and adaptive filtering using a pretrained optical flow network. (iii) **DVF [32]** uses volumetric

4

| Method | Inputs | Vimeo-90K | | UCF101 | | DAVIS | |
|---|---|---|---|---|---|---|---|
| | | PSNR (↑) | SSIM(↑) | PSNR(↑) | SSIM(↑) | PSNR(↑) | SSIM(↑) |
| DAIN [1] | RGB+Depth+Flow | 33.35 | 0.945 | 31.64 | 0.957 | 26.12 | 0.870 |
| QVI [76] | RGB+Flow | 35.15 | 0.971 | 32.89 | 0.970 | 27.17 | **0.874** |
| DVF [32] | RGB | 27.27 | 0.893 | 28.72 | 0.937 | 22.13 | 0.800 |
| SepConv [43] | RGB | 33.60 | 0.944 | 31.97 | 0.943 | 26.21 | 0.857 |
| CAIN [9] | RGB | 33.93 | 0.964 | 32.28 | 0.965 | 26.46 | 0.856 |
| SuperSloMo [22] | RGB | 32.90 | 0.957 | 32.33 | 0.960 | 25.65 | 0.857 |
| BMBC [47] | RGB | 34.76 | 0.965 | 32.61 | 0.955 | 26.42 | 0.868 |
| AdaCoF [27] | RGB | 35.40 | 0.971 | 32.71 | 0.969 | 26.49 | 0.866 |
| FLAVR | RGB | **36.30** | **0.975** | **33.33** | **0.971** | **27.44** | **0.874** |

Table 1: **Comparison with state-of-the-art methods for 2x interpolation** on Vimeo-90K, UCF101, and DAVIS datasets. The upper table includes the methods that use additional input such as optical flows and/or depth maps. The lower table represents the methods using only RGB as input. The first and second best methods are marked in bold and underlined text. Our method consistently outperforms prior works which take only RGB as input, as well as works which additionally require optical flows and/or depth inputs.

sampling to generate the output frame using voxel optical flow predictions. (iv) **SepConv [43]** predicts optimum pairs of spatially varying kernels used to resample from input images in a local neighborhood. (v) **SuperSloMo [22]** performs warping based on flow and visibility maps computed from the input frames. (vi) **CAIN [9]** performs frame interpolation by channel attention and sub-pixel convolutions, and is limited to single frame interpolation, and (vii) **FLAVR** is our proposed approach. To enable fair comparison with our method, we retrain the networks of [1], [76], [22], and [9] on the same datasets (Vimeo, GoPro) as ours. For [32] and [43], we use the pretrained models provided by the authors.

## 5. Results and Analysis

In this section, we compare with the state-of-the-art approaches for video frame interpolation and present analysis and insights into our approach.

### 5.1. Comparison with the state-of-the-art

**Single-Frame Interpolation.** We report the results for single frame interpolation in Table. 1, corresponding to $2\times(k{=}2)$ interpolation from 15 FPS to 30 FPS. We observe that FLAVR outperforms prior methods by a significant margin on Vimeo-90K dataset and sets the *new state-of-the-art* on this dataset with a PSNR value of 36.3 and SSIM value of 0.975. FLAVR is a more generally applicable method and outperforms [22, 32, 43] which assume uniform linear motion between the frames. FLAVR also performs better than [9] which uses a similar end to end architecture to predict output frames, underlining the benefits achieved using an encoder-decoder architecture with spatio-temporal kernels. More importantly, FLAVR also outperforms DAIN [1]

| Method | Inputs | Adobe | | GoPro | |
|---|---|---|---|---|---|
| | | PSNR | SSIM | PSNR | SSIM |
| DAIN [1] | RGB+Depth+Flow | 29.50 | 0.910 | 29 | 0.91 |
| QVI [76] | RGB+Flow | **33.68** | **0.97** | 30.55 | 0.933 |
| DVF [32] | RGB | 28.23 | 0.896 | 21.94 | 0.776 |
| SuperSloMo [22] | RGB | 30.66 | 0.391 | 28.52 | 0.891 |
| FLAVR | RGB | 32.20 | 0.957 | **31.31** | **0.94** |

Table 2: **Comparison with state-of-the-art methods for 8x interpolation** on Adobe and GoPro datasets. FLAVR outperforms all previous work that use only RGB as input. Overall, FLAVR achieves the best performance on GoPro and the second best on Adobe.

and QVI [76] without demanding additional knowledge in the form of bidirectional flow or depth maps.

We test the generalization capability of our method by evaluating the same trained model on UCF101 and DAVIS datasets. These are relatively more challenging for video frame interpolation, containing complex object and human motions from a range of dynamic scenes. Nevertheless, with a PSNR of 33.33 on the UCF101 dataset and 27.44 on the DAVIS dataset, FLAVR clearly delivers better performance compared to all the baselines methods which take RGB images as inputs, and performs on par or better than methods that additionally demand depth or flow maps as inputs.

**Multi-Frame Interpolation.** For multi-frame setting, we report results on $8\times$ ($k{=}8$) interpolation in Table. 2, which corresponds to going from 30 to 240 FPS by generating 7 intermediate frames. Our method yields a PSNR of 31.31 and an SSIM score of 0.94 on the GoPro dataset, which is better than all the prior approaches proposed for frame interpolation. On the Adobe dataset, our method outperforms all methods significantly except QVI, but QVI additionally uses an optical flow estimator which we do
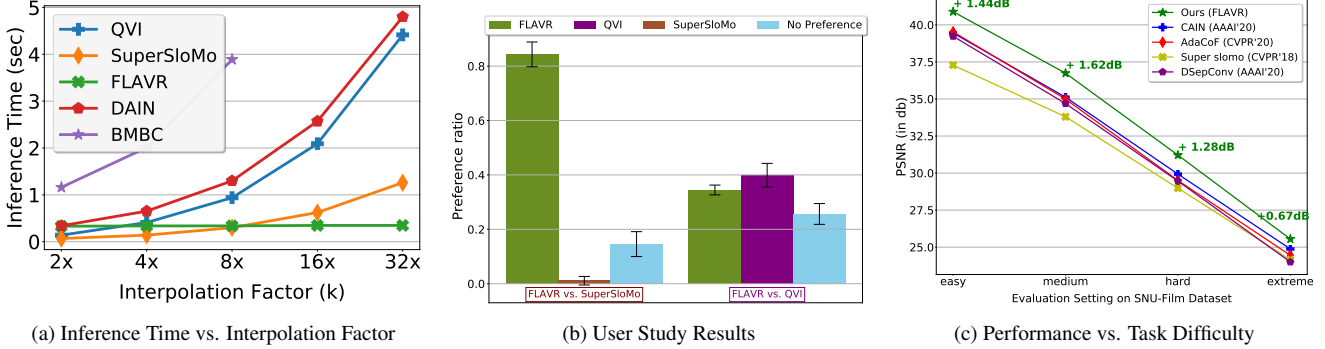
| | (a) Inference Time vs. Interpolation Factor | (b) User Study Results | (c) Performance vs. Task Difficulty |

Figure 3: **Analysis**. (a) Inference time (forward pass w/o IO) of different methods on different interpolation factor. FLAVR has almost no change in inference time due to its design to predict multiple frames per inference. (b) Comparison between FLAVR with Super-SloMo and QVI in a user study on DAVIS. FLAVR significantly outperforms Super-SloMo, and performs comparable to QVI. (c) Comparison between FLAVR with other methods on SNU-Film dataset. FLAVR consistently outperforms all comparing methods across all levels of task difficulty.

not. Similar improvements in performance can also be observed in the case of $4\times$ ($k$=4) interpolation, as shown in Table. 3. These results indicate the effectiveness of the proposed FLAVR architecture even for the task of multi-frame interpolation.

**Middlebury Dataset.** We also evaluate FLAVR on the publicly available test images from Middleburry dataset [53] on the task of single frame interpolation. However, Middleburry has test samples with only two input frames while FLAVR requires 4-frame inputs. In those examples, we simply duplicate them into 4 frames and evaluate with FLAVR. For two frame sequences like *teddy*, duplicating inputs is obviously sub-optimal. On sequences where multi-frame inputs are available, FLAVR outperforms most prior interpolation works like SuperSloMo [22], BMBC [47] and EDSC [5]. Qualitative results for some such sequences are presented in Figure. 8. FLAVR ranks **1st, 4th, 7th** on *backyard, evergreen, basketball* sequences respectively, at the time of this submission. The complete results are available on the public leaderboard.

**Speed vs. Accuracy Trade-off.** One major challenge for realizing the applications of video frame interpolation for real time applications on low resource hardware is to optimize the trade off between faster inference speed and better interpolation quality. Perhaps the most important contribution of our work is to propose an approach that strikes an optimum balance between both these factors by achieving best performance with smallest runtime. As show in Figure. 1, FLAVR offers an improved run time for multi-frame interpolation models. This improvement is possible mainly because we require no overhead in terms of computing optical flow or depth, and predict all the frames in a single forward pass. We also show in Figure. 3a that the inference speed using our method scales gracefully with an increase in the interpolation factor $k$, while most prior methods incur linear growth with $k$. We achieve runtime improvements of

| | SuperSloMo [22] | DAIN [1] | QVI [76] | FLAVR |
|---|---|---|---|---|
| PSNR | 30.8 | 32.49 | <u>36.29</u> | **37.82** |
| SSIM | 0.924 | 0.957 | <u>0.980</u> | **0.983** |

Table 3: **Comparison with state-of-the-art methods for 4x interpolation** on Adobe dataset. FLAVR outperforms QVI by 1.6dB, DAIN by 5.3dB, and Super-SloMo by 7dB in PSNR.

$2.7\times$ and $6.2\times$ for $8\times$ and $16\times$ interpolation respectively with respect to QVI, the current most accurate method, which is possible due to our single shot flow-free prediction and efficient architecture.

**Robustness to Task Difficulty.** We validate the robustness in performance of our method using the SNU-Film dataset [9] consisting of videos with varying difficulty for interpolation depending on the temporal gap between the input frames. The four settings we use are *easy* (120-240 FPS), e.g. predicting 240 FPS video from 120 FPS input, *medium* (60-120 FPS), *hard* (30-60 FPS) and *extreme* (15-30 FPS). In Figure. 3c, we compare the performance of our method with prior works including CAIN [9] and AdaCoF [27], and report better performance than all the methods consistently across all the difficulty settings. Specifically, we see a gain of $1.28dB$ and $1.62dB$ compared to the next best approach [9] in the *hard* and *medium* settings respectively, which are considered challenging for video frame interpolation because of large motions and longer time gaps between the frames. This indicates that our method is more robust to input video frame rates.

**User Study.** Commonly used quantitative metrics like PSNR and SSIM do not strongly correlate with human visual perception to image quality [44, 17]. Therefore, we conduct a large scale user survey on Amazon Mechanical Turk to analyze the performance of our method in comparison to competing methods. Starting from the 90 HD videos from the DAVIS 2017 dataset [50], we generate SloMo videos
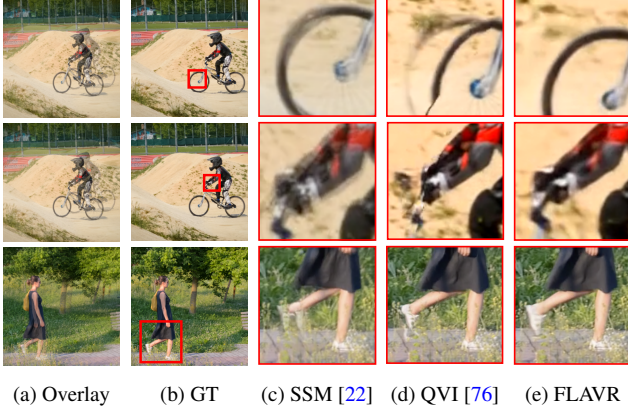
6

| (a) Overlay | (b) GT | (c) SSM [22] | (d) QVI [76] | (e) FLAVR |

Figure 4: **Qualitative comparison with state-of-the-art methods**. We qualitatively compare FLAVR with Super-SloMo (SSM), QVI on a few video sequences on DAVIS.

using $8\times$ interpolation using our method as well as QVI [76] and Super-SloMo [22] for comparison. We explain the details of the survey in Appendix D, and summarize the results in Figure. 3b.

Firstly, when the comparison is between our method against Super-SloMo, users overwhelmingly preferred our videos as the generated videos looked more realistic with minimum artefacts around edges and motion boundaries owing to accurate interpolation. In comparison with QVI, users choose FLAVR in 35% of videos, compared to QVI which was chosen in 40% of the videos, and for 20% of videos the differences came out to be negligible. These results further support our hypothesis that in the interest of real world deployment, optical flow and warping based frame interpolation methods can be substituted with our learning based approach that offers orders of magnitude faster inference (Figure. 3a) with minimal loss in performance.

### 5.2. Ablations

In Table. 4, we present a detailed ablation study of the proposed architecture design in terms of the skip connections, strides and loss functions. We explain each of them in detail next. We conduct all the ablation studies on the Vimeo-90K dataset.

**Backbone Architecture and Input Context** In this work, we propose using 3D convolutions that model space-time relations for improved frame interpolation. To verify this hypothesis, we train a video interpolation network using 2D convolutions instead, and present the results in Table. 4a. While training 2D Resnet, we concatenate RGB channels of the input before feeding into the network. We observe that the *R2D-18-2I* baseline, which uses a 2D ResNet-18 encoder decoder with 2 input frames ($C = 1$) performs worse than *2D-R18-4I* baseline, which uses 4 input frames ($C = 2$) justifying the need for a larger input context. Next, our proposed architecture *3D-R18-4I* which uses 3D convolutions
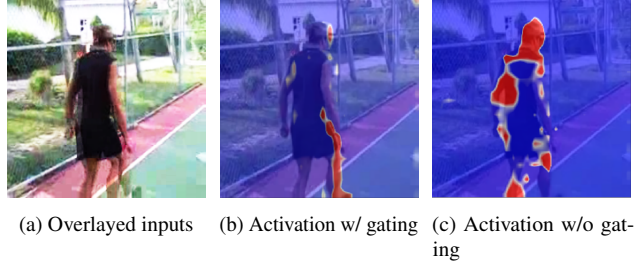


| (a) Overlayed inputs | (b) Activation w/ gating | (c) Activation w/o gating |

Figure 5: **Visualization of attention weighted feature maps** (a) The overlayed input frames. (b) The feature map of the channel with the highest attention weight in the network with feature gating. (c) The same feature map without using the gating module. We observe higher activation (red) in (b) along the motion boundaries. Best viewed in color.

along with 4 inputs, clearly outperforms both these baselines by $1.3$ and $2.3$dB, respectively. This indicates the importance of temporal modeling for the task. We also experiment with longer inputs (e.g. $C = 3$ input frames, *R3D-18-6I*), but observe no gain while requiring more computation. This can be explained as the outer frames generally contain very little relevant detail for interpolation due to large motion and shot changes. Therefore we settle on using 4 input frames, or $C = 2$, in all other experiments. Finally, we did not observe improvement while using a deeper backbone like R3D-34, mostly because of overfitting issue due to limited training data.

**Choice of Fusion** Table. 4b compares and reports the different choices for the skip connection (in Figure. 2) used for combining features across encoder and corresponding decoder. *No fusion* corresponds to having no skip connection between the layers of the encoder and decoder. While *fusion - add* corresponds to adding the features from the encoder to the decoder, *fusion - concat* refers to concatenating the corresponding feature maps along the channels. We find that using some kind of feature transfer across encoder and decoder is essential, than having *No fusion* (PSNR of 36.11 vs. 35.1), as the complementary information learnt in the low level and high level features needs to be aggregated for accurate interpolation. We settle on using *fusion - concat* in our final model as it gives better performance than *fusion - add*.

**Temporal Striding** Striding or pooling in CNNs are known to remove lot of fine level details in images, which are essential for generative tasks like frame interpolation. We verify this with experiments using $2\times(1/2\times)$ and $4\times(1/4\times)$ temporal striding in the encoder(decoder), and observe from Table. 4c that the performance decreases from 36.3 to 35.2 with larger temporal striding. We conclude that temporal striding hurts, and use a temporal stride of 1 in all the 3D convolution layers.

**Channel Gating** We visualize the role of channel gating module in the network in Figure. 5. We show the over-

| Model | PSNR | SSIM |
|---|---|---|
| R2D-18-2I | 33.98 | 0.966 |
| R2D-18-4I | 34.97 | 0.967 |
| R3D-18-4I | **36.3** | **0.975** |
| R3D-18-6I | 35.47 | 0.971 |
| R3D-34-4I | 35.85 | 0.972 |

(a) Effect of encoder arch.

| Model | PSNR | SSIM |
|---|---|---|
| No fusion | 35.1 | 0.9713 |
| fusion - add | 35.7 | 0.9737 |
| fusion - concat | **36.3** | **0.975** |

(b) Type of feature fusion from encoder and decoder

| Model | PSNR | SSIM |
|---|---|---|
| w/o stride | **36.3** | **0.975** |
| w/ 2x stride | 35.4 | 0.961 |
| w/ 4x stride | 35.21 | 0.96 |

(c) Effect of larger temporal striding

| Model | PSNR | SSIM |
|---|---|---|
| L1 Loss | **36.3** | **0.975** |
| L2 Loss | 35.3 | 0.965 |
| Huber Loss | 35.3 | 0.964 |
| L1+VGG Loss | 35.91 | 0.962 |

(d) Effect of loss function

Table 4: **Ablation results**. Ablation results for FLAVR architecture on (a) different backbones, (b) fusion methods, (c) temporal striding, and (d) loss functions.

| Method | pretrained on | Arch. | UCF101 | HMDB51 |
|---|---|---|---|---|
| Random Init. | - | R3D-18 | 50.02 | 19.00 |
| Supervised | Kinetics-400 | R3D-18 | 87.70 | 59.10 |
| Contrastive [14] | Kinetics-400 | R3D-18 | 68.20 | 34.50 |
| Video-GAN [62] | UCF101 | Custom | 52.10 | - |
| LMD[33] | NTU RGB | Custom | 53.00 | - |
| DVF [32] | UCF101 | Custom | 52.40 | - |
| FLAVR | Vimeo-90K | R3D-18 | 63.10 | 23.48 |

Table 5: **FLAVR for action recognition**. Results of using FLAVR as a self-supervised pretext task for action recognition on UCF101 and HMDB51. The upper table includes the baselines and the state-of-the-art contrastive-based self-supervised learning method, DPC [14]. The lower table presents self-supervised approaches with pixel-level pretext tasks. FLAVR significantly outperforms the training from scratch baseline and other self-supervised approaches that use pixel-level pretext tasks.

lapped input frames in Figure. 5a to highlight the parts which have motion. In Figure. 5b and Figure. 5c, we plot the feature maps corresponding to the channel dimension with the largest activation while using and without using the feature gating respectively. We observe that the network trained with spatio-temporal gating (Figure. 5b) learns to focus on parts of input with visible motions (high activations in red), thus resulting in confident predictions of the interpolated motion estimates compared to Figure. 5c. In fact, training without spatio-temporal gating results in a drop in PSNR value from 36.3 to 36.1, further validating the utility of having the gating module.

**Loss Function** Many previous works [43] have studied the effect of using purely pixel loss vs. perception based losses [23]. Using only L1 or L2 loss would improve on the PSNR metric, but would cause blur in predictions. On the other hand, adding VGG based perception loss would result in sharper images visually. We observe from Table. 4d that we did not improve upon the PSNR or SSIM metric by using any additional loss functions like VGG loss or Huber loss, apart from just L1 loss which also resulted in visually sharper images in our case.

## 6. FLAVR as a self-supervised pretext task

**Action Recognition** We evaluate the semantic properties of the internal representations learned by FLAVR by reusing its trained encoder on a downstream action recognition task. We remove the decoder and attach a classification layer to the network. The whole network is then finetuned end to end on UCF101 and HMDB51 datasets. In order to isolate the benefits arising from pretraining the encoder on video interpolation task, we train a 3D resnet (R3D) baseline completely from scratch and observe from Table. 5 that FLAVR, which is pretrained on Vimeo-90K dataset on frame interpolation task clearly outperforms random initialization baseline by $13.08\%$ on UCF-101 and $4.48\%$ on HMDB-51. FLAVR also significantly outperforms prior self-supervised methods on video which use low level pretext tasks like Video-GAN [62] and flow descriptors [33] indicating that frame interpolation can learn useful motion representations. Finally, FLAVR also achieves better accuracy than pretraining using DVF [32] which underlines that our particular method for video frame interpolation invariably benefits downstream action recognition more than voxel flow.

**Optical Flow Estimation** It is known that successful frame interpolation intrinsically depends on reliable optical flow estimation [70]. We investigate this hypothesis on our method by finetuning our trained network for optical flow estimation on MPI Sintel [3] validation set, and report the corresponding EPE (end point error) in Table. 6. Finetuning using FLAVR achieves much lower EPE compared with random initialization using the same backbone architecture, proving that our model learns useful flow features. We note that we do not aim to outperform more complex, flow-dedicated architectures[19, 29] but only aim to understand if we can learn useful flow features using a simple architecture like ours by pre-training on frame interpolation.

**Motion Magnification** Motion magnification [68, 46] is a complementary problem to frame interpolation, in which the task is to magnify the subtle motions from the input video. Using FLAVR as pre-training, we fine-tune a motion magnification network for a fixed magnification factor of 10. FLAVR achieves an SSIM of 0.801 on the synthetic CoCo-Synth dataset [46] with a simple architecture. More

| | FlowNet2.0 [19] | SelFlow [29] | Random Init. | Finetune on FLAVR |
|---|---|---|---|---|
| MPI-Clean | 2.02 | 1.68 | 12.15 | 6.14 |
| MPI-Final | 3.14 | 1.77 | 11.46 | 6.32 |

Table 6: **FLAVR for optical flow estimation**. Comparison between FLAVR and other approaches for optical flow prediction on MPI dataset [3]. FLAVR pre-training helps to reduce the EPE by 6.01 and 5.14 on MPI-Clean and MPI-Final, respectively compared with training from scratch on the same backbone architecture.

details and qualitative results are presented in Appendix A and project video.

## 7. Conclusion

In this paper, we present FLAVR for flow-free and end-to-end video frame interpolation. FLAVR uses 3D convolutions to explicitly model the spatio-temporal relations between the input frames which improves interpolation accuracy under complex motion profiles across various input frame rates. FLAVR offers best trade-off in terms of inference speed vs. interpolation accuracy, while demonstrating significantly better accuracy compared to many existing approaches. We demonstrate that the representations learned by FLAVR are useful for various downstream tasks such as action recognition, optical flow estimation, and motion magnification. We will be publicly releasing all the code and models used in this work upon publication.

## Acknowledgements

## References

[1] 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 Conference on Computer Vision and Pattern Recognition*, pages 3703–3712, 2019. 1, 2, 3, 4, 5, 6

[2] Wenbo Bao, Wei-Sheng Lai, Xiaoyun Zhang, Zhiyong Gao, and Ming-Hsuan Yang. Memc-net: Motion estimation and motion compensation driven neural network for video interpolation and enhancement. *IEEE transactions on pattern analysis and machine intelligence*, 2019. 1, 2

[3] D. J. Butler, J. Wulff, G. B. Stanley, and M. J. Black. A naturalistic open source movie for optical flow evaluation. In A. Fitzgibbon et al. (Eds.), editor, *European Conf. on Computer Vision (ECCV)*, Part IV, LNCS 7577, pages 611–625. Springer-Verlag, Oct. 2012. 8, 9

[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. 3, 13

[5] Xianhang Cheng and Zhenzhong Chen. Multiple video frame interpolation via enhanced deformable separable convolution. *arXiv preprint arXiv:2006.08070*, 2020. 6, 14

[6] Xianhang Cheng and Zhenzhong Chen. Video frame interpolation via deformable separable convolution. In *AAAI*, pages 10607–10614, 2020. 1, 2

[7] Zhixiang Chi, Rasoul Mohammadi Nasiri, Zheng , Juwei Lu, Jin Tang, and Konstantinos N Plataniotis. All at once: Temporally adaptive multi-frame interpolation with advanced motion modeling. *arXiv preprint arXiv:2007.11762*, 2020. 1, 2

[8] Myungsub Choi, Janghoon Choi, Sungyong Baik, Tae Hyun Kim, and Kyoung Mu Lee. Scene-adaptive video frame interpolation via meta-learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9444–9453, 2020. 2

[9] Myungsub Choi, Heewon Kim, Bohyung Han, Ning Xu, and Kyoung Mu Lee. Channel attention is all you need for video frame interpolation. In *AAAI*, pages 10663–10671, 2020. 1, 2, 4, 5, 6

[10] Carl Doersch, Abhinav Gupta, and Alexei A Efros. Unsupervised visual representation learning by context prediction. In *Proceedings of the IEEE international conference on computer vision*, pages 1422–1430, 2015. 2

[11] Christoph Feichtenhofer, Haoqi Fan, Jitendra Malik, and Kaiming He. Slowfast networks for video recognition. In *Proceedings of the IEEE international conference on computer vision*, pages 6202–6211, 2019. 3

[12] Basura Fernando, Hakan Bilen, Efstratios Gavves, and Stephen Gould. Self-supervised video representation learning with odd-one-out networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3636–3645, 2017. 2

[13] Tobias Gurdan, Martin R Oswald, Daniel Gurdan, and Daniel Cremers. Spatial and temporal interpolation of multi-view image sequences. In *German Conference on Pattern Recognition*, pages 305–316. Springer, 2014. 2

[14] Tengda Han, Weidi Xie, and Andrew Zisserman. Video representation learning by dense predictive coding. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*, 2019. 2, 8

[15] Tengda Han, Weidi Xie, and Andrew Zisserman. Memory-augmented dense predictive coding for video representation learning. *arXiv preprint arXiv:2008.01065*, 2020. 2

[16] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9729–9738, 2020. 2

[17] Alain Hore and Djemel Ziou. Image quality metrics: Psnr vs. ssim. In *2010 20th international conference on pattern recognition*, pages 2366–2369. IEEE, 2010. 6

[18] Jie Hu, Li Shen, and Gang Sun. Squeeze-and-excitation networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7132–7141, 2018. 4

[19] Eddy Ilg, Nikolaus Mayer, Tonmoy Saikia, Margret Keuper, Alexey Dosovitskiy, and Thomas Brox. Flownet 2.0:

Evolution of optical flow estimation with deep networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2462–2470, 2017. 8, 9

[20] Allan Jabri, Andrew Owens, and Alexei A Efros. Space-time correspondence as a contrastive random walk. *arXiv preprint arXiv:2006.14613*, 2020. 2

[21] Dinesh Jayaraman and Kristen Grauman. Slow and steady feature analysis: higher order temporal coherence in video. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3852–3861, 2016. 2

[22] Huaizu Jiang, Deqing Sun, Varun Jampani, Ming-Hsuan Yang, Erik Learned-Miller, and Jan Kautz. Super slomo: High quality estimation of multiple intermediate frames for video interpolation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 9000–9008, 2018. 1, 2, 3, 5, 6, 7, 12, 13, 14

[23] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *European conference on computer vision*, pages 694–711. Springer, 2016. 8

[24] Soomro Khurram, Zamir Amir, and Shah Mubarak. UCF101: A dataset of 101 human action classes from videos in the wild. *CRCV-TR-12-01*, 2012. 4, 12

[25] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 4

[26] Hildegard Kuehne, Hueihan Jhuang, Estíbaliz Garrote, Tomaso Poggio, and Thomas Serre. Hmdb: a large video database for human motion recognition. In *2011 International Conference on Computer Vision*, pages 2556–2563. IEEE, 2011. 12

[27] Hyeongmin Lee, Taeoh Kim, Tae-young Chung, Daehyun Pak, Yuseok Ban, and Sangyoun Lee. Adacof: Adaptive collaboration of flows for video frame interpolation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5316–5325, 2020. 1, 5, 6

[28] Hsin-Ying Lee, Jia-Bin Huang, Maneesh Singh, and Ming-Hsuan Yang. Unsupervised representation learning by sorting sequences. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 667–676, 2017. 2

[29] Pengpeng Liu, Michael Lyu, Irwin King, and Jia Xu. Selflow: Self-supervised learning of optical flow. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4571–4580, 2019. 8, 9

[30] Yihao Liu, Liangbin Xie, Li Siyao, Wenxiu Sun, Yu Qiao, and Chao Dong. Enhanced quadratic video interpolation. In *European Conference on Computer Vision*, pages 41–56. Springer, 2020. 2

[31] Yu-Lun Liu, Yi-Tung Liao, Yen-Yu Lin, and Yung-Yu Chuang. Deep video frame interpolation using cyclic frame generation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 8794–8802, 2019. 2

[32] Ziwei Liu, Raymond A Yeh, Xiaoou Tang, Yiming Liu, and Aseem Agarwala. Video frame synthesis using deep voxel flow. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 4463–4471, 2017. 1, 2, 4, 5, 8

[33] Zelun Luo, Boya Peng, De-An Huang, Alexandre Alahi, and Li Fei-Fei. Unsupervised learning of long-term motion dynamics for videos. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2203–2212, 2017. 8

[34] Dhruv Mahajan, Fu-Chung Huang, Wojciech Matusik, Ravi Ramamoorthi, and Peter Belhumeur. Moving gradients: a path-based method for plausible image interpolation. *ACM Transactions on Graphics (TOG)*, 28(3):1–11, 2009. 2

[35] Simone Meyer, Abdelaziz Djelouah, Brian McWilliams, Alexander Sorkine-Hornung, Markus Gross, and Christopher Schroers. Phasenet for video frame interpolation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 498–507, 2018. 1, 2

[36] Simone Meyer, Oliver Wang, Henning Zimmer, Max Grosse, and Alexander Sorkine-Hornung. Phase-based frame interpolation for video. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1410–1418, 2015. 2

[37] Antoine Miech, Ivan Laptev, and Josef Sivic. Learnable pooling with context gating for video classification. *arXiv preprint arXiv:1706.06905*, 2017. 4

[38] Ishan Misra, C Lawrence Zitnick, and Martial Hebert. Shuffle and learn: unsupervised learning using temporal order verification. In *European Conference on Computer Vision*, pages 527–544. Springer, 2016. 2

[39] Seungjun Nah, Tae Hyun Kim, and Kyoung Mu Lee. Deep multi-scale convolutional neural network for dynamic scene deblurring. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3883–3891, 2017. 4

[40] Simon Niklaus and Feng Liu. Context-aware synthesis for video frame interpolation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1701–1710, 2018. 1, 2

[41] 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. 1, 2

[42] Simon Niklaus, Long Mai, and Feng Liu. Video frame interpolation via adaptive convolution. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 670–679, 2017. 1, 2

[43] 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. 1, 2, 5, 8

[44] Jim Nilsson and Tomas Akenine-Möller. Understanding ssim. *arXiv preprint arXiv:2006.13846*, 2020. 4, 6

[45] Augustus Odena, Vincent Dumoulin, and Chris Olah. Deconvolution and checkerboard artifacts. *Distill*, 2016. 4

[46] Tae-Hyun Oh, Ronnachai Jaroensri, Changil Kim, Mohamed Elgharib, Fr'edo Durand, William T Freeman, and Wojciech Matusik. Learning-based video motion magnification. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 633–648, 2018. 8, 12

[47] Junheum Park, Keunsoo Ko, Chul Lee, and Chang-Su Kim. Bmbc: Bilateral motion estimation with bilateral cost volume

for video interpolation. *arXiv preprint arXiv:2007.12622*, 2020. 1, 5, 6, 14

[48] Deepak Pathak, Ross Girshick, Piotr Dollár, Trevor Darrell, and Bharath Hariharan. Learning features by watching objects move. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2701–2710, 2017. 2

[49] Tomer Peleg, Pablo Szekely, Doron Sabo, and Omry Sendik. Im-net for high resolution video frame interpolation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2398–2407, 2019. 1, 2

[50] Federico Perazzi, Jordi Pont-Tuset, Brian McWilliams, Luc Van Gool, Markus Gross, and Alexander Sorkine-Hornung. A benchmark dataset and evaluation methodology for video object segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 724–732, 2016. 4, 6

[51] 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 International Conference on Computer Vision*, pages 892–900, 2019. 2

[52] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pages 234–241. Springer, 2015. 3

[53] Daniel Scharstein, Heiko Hirschmüller, York Kitajima, Greg Krathwohl, Nera Nešić, Xi Wang, and Porter Westling. High-resolution stereo datasets with subpixel-accurate ground truth. In *German conference on pattern recognition*, pages 31–42. Springer, 2014. 6

[54] O. Shahar, A. Faktor, and M. Irani. Space-time super-resolution from a single video. In *CVPR 2011*, pages 3353–3360, 2011. 2

[55] Zhihao Shi, Xiaohong Liu, Kangdi Shi, Linhui Dai, and Jun Chen. Video interpolation via generalized deformable convolution. *arXiv preprint arXiv:2008.10680*, 2020. 2

[56] Mihoko Shimano, Takahiro Okabe, Imari Sato, and Yoichi Sato. Video temporal super-resolution based on self-similarity. In *Asian Conference on Computer Vision*, pages 93–106. Springer, 2010. 2

[57] Zheng Shou, Jonathan Chan, Alireza Zareian, Kazuyuki Miyazawa, and Shih-Fu Chang. CDC: Convolutional-de-convolutional networks for precise temporal action localization in untrimmed videos. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017. 3

[58] Shuochen Su, Mauricio Delbracio, Jue Wang, Guillermo Sapiro, Wolfgang Heidrich, and Oliver Wang. Deep video deblurring for hand-held cameras. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1279–1288, 2017. 2, 4

[59] 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 conference on computer vision and pattern recognition*, pages 8934–8943, 2018. 2

[60] Du Tran, Lubomir Bourdev, Rob Fergus, Lorenzo Torresani, and Manohar Paluri. Learning spatiotemporal features with 3d convolutional networks. In *Proceedings of the IEEE international conference on computer vision*, pages 4489–4497, 2015. 3

[61] Du Tran, Heng Wang, Lorenzo Torresani, Jamie Ray, Yann LeCun, and Manohar Paluri. A closer look at spatiotemporal convolutions for action recognition. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pages 6450–6459, 2018. 3

[62] Carl Vondrick, Hamed Pirsiavash, and Antonio Torralba. Generating videos with scene dynamics. In *Advances in neural information processing systems*, pages 613–621, 2016. 8

[63] Carl Vondrick, Abhinav Shrivastava, Alireza Fathi, Sergio Guadarrama, and Kevin Murphy. Tracking emerges by colorizing videos. In *Proceedings of the European conference on computer vision (ECCV)*, pages 391–408, 2018. 2

[64] Neal Wadhwa, Michael Rubinstein, Frédo Durand, and William T. Freeman. Phase-based video motion processing. *ACM Trans. Graph. (Proceedings SIGGRAPH 2013)*, 32(4), 2013. 12

[65] Xiaolong Wang and Abhinav Gupta. Unsupervised learning of visual representations using videos. In *Proceedings of the IEEE international conference on computer vision*, pages 2794–2802, 2015. 2

[66] Xiaolong Wang, Allan Jabri, and Alexei A Efros. Learning correspondence from the cycle-consistency of time. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2566–2576, 2019. 2

[67] Donglai Wei, Joseph J Lim, Andrew Zisserman, and William T Freeman. Learning and using the arrow of time. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8052–8060, 2018. 2

[68] Hao-Yu Wu, Michael Rubinstein, Eugene Shih, John Guttag, Frédo Durand, and William Freeman. Eulerian video magnification for revealing subtle changes in the world. *ACM transactions on graphics (TOG)*, 31(4):1–8, 2012. 8, 12

[69] Junru Wu, Xiang Yu, Ding Liu, Manmohan Chandraker, and Zhangyang Wang. David: Dual-attentional video deblurring. In *The IEEE Winter Conference on Applications of Computer Vision*, pages 2376–2385, 2020. 2

[70] Jonas Wulff and Michael J Black. Temporal interpolation as an unsupervised pretraining task for optical flow estimation. In *German Conference on Pattern Recognition*, pages 567–582. Springer, 2018. 8

[71] 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. *CoRR*, abs/2002.11616, 2020. 2

[72] Saining Xie, Chen Sun, Jonathan Huang, Zhuowen Tu, and Kevin Murphy. Rethinking spatiotemporal feature learning: Speed-accuracy trade-offs in video classification. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 305–321, 2018. 3, 4

[73] Dejing Xu, Jun Xiao, Zhou Zhao, Jian Shao, Di Xie, and Yueting Zhuang. Self-supervised spatiotemporal learning via video clip order prediction. In *Proceedings of the IEEE*

*Conference on Computer Vision and Pattern Recognition*, pages 10334–10343, 2019. 2

[74] Huijuan Xu, Abir Das, and Kate Saenko. R-c3d: Region convolutional 3d network for temporal activity detection. In *Proceedings of the International Conference on Computer Vision (ICCV)*, 2017. 3

[75] Jun Xu, Tao Mei, Ting Yao, and Yong Rui. Msr-vtt: A large video description dataset for bridging video and language. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016. 3

[76] Xiangyu Xu, Li Siyao, Wenxiu Sun, Qian Yin, and Ming-Hsuan Yang. Quadratic video interpolation. In *Advances in Neural Information Processing Systems*, pages 1647–1656, 2019. 1, 2, 4, 5, 6, 7, 12

[77] Tianfan Xue, Baian Chen, Jiajun Wu, Donglai Wei, and William T Freeman. Video enhancement with task-oriented flow. *International Journal of Computer Vision*, 127(8):1106–1125, 2019. 1, 2, 4

[78] Zhefei Yu, Houqiang Li, Zhangyang Wang, Zeng Hu, and Chang Wen Chen. Multi-level video frame interpolation: Exploiting the interaction among different levels. *IEEE Transactions on Circuits and Systems for Video Technology*, 23(7):1235–1248, 2013. 2

[79] Liangzhe Yuan, Yibo Chen, Hantian Liu, Tao Kong, and Jianbo Shi. Zoom-in-to-check: Boosting video interpolation via instance-level discrimination. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 12183–12191, 2019. 2

[80] Haoxian Zhang, Yang Zhao, and Ronggang Wang. A flexible recurrent residual pyramid network for video frame interpolation. ICCV, 2019. 2

[81] Liad Pollak Zuckerman, Shai Bagon, Eyal Naor, George Pisha, and Michal Irani. Across scales & across dimensions: Temporal super-resolution using deep internal learning. *arXiv preprint arXiv:2003.08872*, 2020. 2

# Appendix

## A. Motion magnification

Motion Magnification [46, 68, 64] deals with magnifying subtle yet important motions from videos, which are often imperceptible by human eyes. From [46], we define motion magnification as follows. For an Image $I(\mathbf{x}, t) = f(\mathbf{x} + \delta(\mathbf{x}, t))$, the goal of motion magnification is to generate an output image $\tilde{I}(\mathbf{x}, t)$ such that

$$\tilde{I}(\mathbf{x}, t) = f(\mathbf{x} + (1 + \alpha)\delta(\mathbf{x}, t)) \qquad (2)$$

for a magnification factor $\alpha$. For frame interpolation, $\alpha < 1$, since we are interested in what happens between two frames while for motion magnification, $\alpha > 1$, since we look to extrapolate existing motions beyond visible regime. While prior works [46, 68, 64] used custom architectures along with various post processing filters for this task, we offer a complementary perspective and look into how much a simple architecture like FLAVR pretrained on frame interpolation

helps motion magnification. For this purpose, we use the synthetic dataset *CoCo-Synth* [46] to perform the training. We train the network for a fixed magnification factor of 10 ($\alpha = 10$). On this dataset, pretraining on FLAVR improved the SSIM values on a held-out validation set from $0.732$ to $0.801$. We provide sample videos after magnification and compare it with phase based approach [68] in our project video. We emphasize that we do not apply any post processing such as temporal or spatial filters for removing noise on the outputs. The videos are generated directly as an output of the FLAVR architecture pretrained on frame interpolation, and finetuned for motion magnification.

## B. Experiment setting for action recognition

For downstream experiments on action recognition, we use the train and validation split 1 of UCF101 [24] and HMDB51 [26]. We remove the decoder from the architecture and use the pretrained encoder along with a classifier (a global average pooling, a fully-connected layer, and a softmax) for training on downstream actions and add a temporal stride of 4. For UCF101, we use an input size of $32 \times 3 \times 224 \times 224$ and for HMDB51 we use an input size of $16 \times 3 \times 224 \times 224$ with a batch size of 16. The networks are fine-tuned using SGD with batch norm with a learning rate of $0.02$ for 40 epochs. During inference, we sample 10 consecutive overlapping clips of length 32 from the test video and average predictions over all the clips.

## C. Experiment setting for optical flow estimation

One crucial point to consider in downstream training on optical flow is that the flow networks generally take only two input frames which is considered too short for 3D CNN. Nevertheless, to examine the effectiveness of features learnt using frame interpolation for optical flow, we use the same encoder and decoder, and initialize the last prediction layer to output two channels instead (corresponding to x and y values of flow at each pixel). Since the interpolation network was trained to take 4-frame inputs, we apply copy padding to the inputs, e.g. repeating each input frame 2 times. We use an EPE (end point error) loss and train our network for 200 epochs. We report numbers using 5-fold cross validation over the MPI-Sintel clean and final subsets. Qualitative results are provided in the project video.

## D. User study

We carry the user study on the Amazon Mechanical Turk (AMT) platform. We select two representative works that belong to two broad families that perform linear (Super-SloMo [22]) and quadratic (QVI [76]) warping for multi-frame interpolation. Then, we compare each video generated by FLAVR with videos generated using each of SuperSloMo

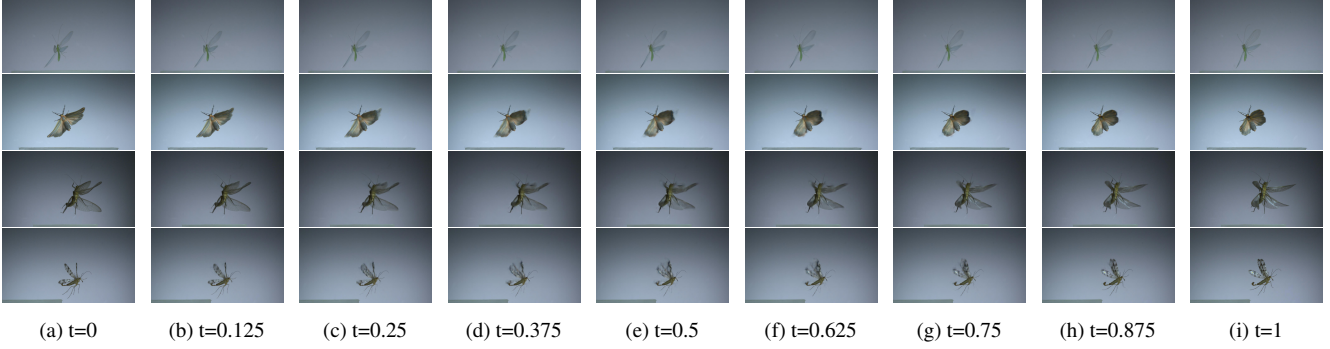|     |     |     |     |     |     |     |     |     |
| (a) t=0 | (b) t=0.125 | (c) t=0.25 | (d) t=0.375 | (e) t=0.5 | (f) t=0.625 | (g) t=0.75 | (h) t=0.875 | (i) t=1 |

Figure 6: Qualitative Results for $8\times$ video frame interpolation on Insect Motion Videos. Frame at $t = 0$ and $t = 1$ are given as inputs to the network to predict the remaining 7 intermediate frames. Original Videos acquired from `AntLab` Youtube Channel.
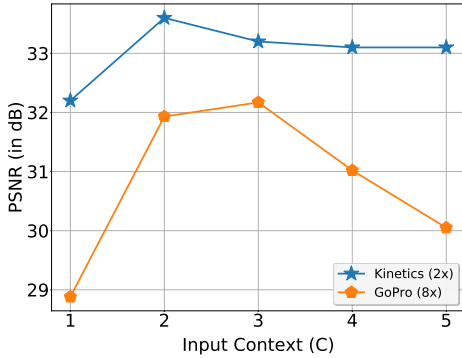


Figure 7: **Effect of Input Context** Comparison of the effect of input context, $C$, for video frame interpolation. For $2\times$ interpolation, we observed that a value of $C = 2$ which corresponds to using 4 input frames, 2 each from the past and future, gives best results. Beyond $C = 3$, we observe no further improvements. For $8\times$ interpolation, a value of $C = 3$ gave the best accuracy.

and QVI separately. For this purpose we use *all* 90 HD videos from the DAVIS dataset, generate $8\times$ interpolated videos and place the two interpolated videos one beside the other and randomly shuffle the order of videos. We then show each pair of videos to 6 AMT workers and ask them to choose which video, right or left, looked more realistic. The method preferred by more users is chosen as a winner for that particular video. In case of tie, that is if each method is chosen by 3 users, we place the video under "no preference" category. With this setting, we find that users overwhelmingly chose our videos in preference against SuperSloMo [22]. More details are provided in section 5.1 of the main paper.

## E. Benchmarking inference time

The inference time benchmarking was performed using an NVIDIA-2080Ti GPU with 12GB memory. The calculated time only includes forward pass excluding the data pre-processing time and CPU/GPU transfer. The results were obtained by averaging over 100 samples from Adobe-240FPS dataset using $512\times512$ crop size. For multi-frame

interpolation, the time required is calculated as the aggregate time required for interpolating all the frames. Non-blocking CUDA operations as well as GPU warm start time were accounted for during inference time computation.

## F. Effect of input context

In Figure. 7, we present a detailed ablation about the effect of input context ($C$) on the performance of interpolation. Since most existing datasets are constrained to only consider 2 ($C=1$) or 4 ($C=2$) input frames, we use raw videos from Kinetics-400 dataset [4] to carry the ablation for $2\times$ interpolation and videos from GoPro dataset for $8\times$ interpolation. We note that Kinetics-400 contains videos with very large motion and intensity variations, abrupt shot changes, and low quality frame due to compression artifacts which makes it not ideal for the task of video frame interpolation. Nevertheless, we use it to perform ablations into the effect of input context. On kinetics dataset, we use a cropsize of 112 and batch size of 32, and for the GoPro dataset, we use a crop size of 256 and batch size of 12 for the purpose of this ablation.

From Figure. 7, we observe that for both $2\times$ and $8\times$ interpolations, using two input frames ($C = 1$), one each from past and present is sub-optimal, as it fails to accurately reason about complex motion profiles and occlusions. Furthermore, for $2\times$ interpolation on Kinetics-400, we found that a value of $C = 2$ gave the best result, and beyond that the performance saturates. This is because the outer frame generally contain less useful information for interpolation and in some cases might contain significant scene shifts which hurts the interpolation accuracy. In the case of $8\times$ interpolation on GoPro video with 240 FPS, the time gap between the frames is tinier, so we find that a value of $C = 3$ performs the best, while any larger value of $C$ hurts the accuracy.
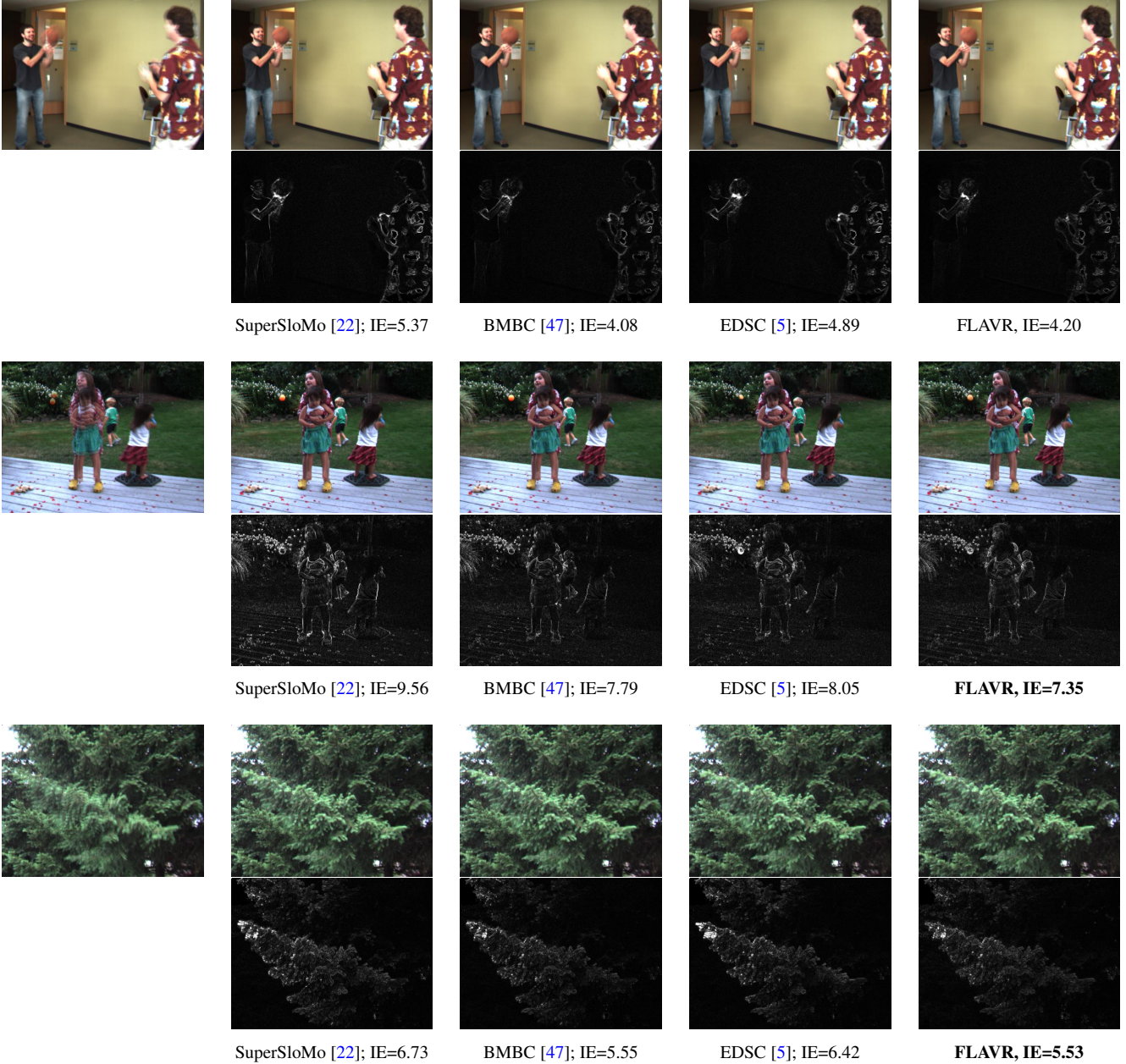
13

Figure 8: Interpolation results for $2\times$ interpolation on Middleburry test set. The leftmost images in each row repesent the overlayed inputs. The first row in each set represents the interpolated frame, while the second row shows the error maps with respect to the ground truth. IE shows the interpolation error of the method. The interpolation errors for all the baselines are reported on the official dashboard.

## G. Training details

We train the $2\times$ interpolation network on Vimeo-90K dataset and $4\times$ and $8\times$ interpolation networks on the GoPro dataset and use the official train and validation splits with the sampling strategy explained in Section 3 of the paper. We use a crop size of $256\times256$ and $512\times512$ for Vimeo-90K and GoPro datasets, respectively. We employ random frame order reversal and random horizontal flipping as augmentation strategies on both the datasets. We use use an initial learning rate of $2 \times 10^{-4}$ and divide the learning rate by 2 whenever the training plateaus. We train the $2\times$ interpolation network for 200 epochs, while $4\times$ and $8\times$ interpolation network were trained for 120 epochs. We use 8 GPUs and a mini-batch of 32 to train each model, and training is completed in about 36 hours for $2\times$ and 22 hours for $8\times$ interpolation networks. To enable fair comparison, all the baseline networks were also trained with exactly similar configurations on the same datasets.

## H. Qualitative Results

We show additional qualitative results by applying frame interpolation technique on insect motion videos. We believe that this application is of immense use for closer inspection of biological properties from videos. We obtain videos from AntLab Youtube channel[2] that have insect takeoff and flying captured at very high FPS. We down-sample the frame rate to 15FPS and apply our interpolation network to recover videos of higher frame rate. We apply our $8\times$ model once to obtain videos of 120FPS. The images are shown in Figure. 6. Complete videos are available in our project video.

---

[2] https://www.youtube.com/user/adrianalansmith