# AIM 2019 Challenge on Video Temporal Super-Resolution: Methods and Results

Seungjun Nah       Sanghyun Son       Radu Timofte       Kyoung Mu Lee       Li Siyao       Ze Pan
Xiangyu Xu       Wenxiu Sun       Myungsub Choi       Heewon Kim       Bohyung Han
Ning Xu       Bumjun Park       Songhyun Yu       Sangmin Kim       Jechang Jeong       Wang Shen
Wenbo Bao       Guangtao Zhai       Li Chen       Zhiyong Gao       Guannan Chen       Yunhua Lu
Ran Duan       Tong Liu       Lijie Zhang       Woonsung Park       Munchurl Kim       George Pisha
Eyal Naor               Lior Aloni

## Abstract

*Videos contain various types and strengths of motions that may look unnaturally discontinuous in time when the recorded frame rate is low. This paper reviews the first AIM challenge on video temporal super-resolution (frame interpolation) with a focus on the proposed solutions and results. From low-frame-rate (15 fps) video sequences, the challenge participants are asked to submit higher-frame-rate (60 fps) video sequences by estimating temporally intermediate frames. We employ the REDS_VTSR dataset derived from diverse videos captured in a hand-held camera for training and evaluation purposes. The competition had 62 registered participants, and a total of 8 teams competed in the final testing phase. The challenge winning methods achieve the state-of-the-art in video temporal super-resolution.*

## 1. Introduction

Video frame interpolation is a classical computer vision task increasing the frame rate of videos. Due to the low-light condition, limited performance of mobile processors, sensors and storage, etc., videos that are recorded in low-frame-rate sometimes exhibit unnatural temporal discontinuity. Video frame interpolation aims to generate the missing views between the recorded frames. For such purpose, video frame interpolation has been adopted in broad fields such as frame rate up-conversion [6, 19, 23, 7, 22] and video coding [14]. Other applications include video slow-motion effect [20] and novel view synthesis [45, 9].

As machine learning has been a recent trend in computer vision, deep neural networks are achieving great successes

in video frame interpolation [26, 25, 34, 35, 33, 13, 52, 20, 27, 37, 2]. However, there has been no standard dataset for video frame interpolation, and most of the existing methods are trained from different datasets. Often, the videos are self-collected from Flickr [34] or YouTube [35, 33, 20, 37]. Other methods propose their own datasets [52, 3] or use exiting datasets introduced for other purposes such as autonomous driving [10], action recognition [44], deblurring [30, 47], and video segmentation [38, 39].

Thus, in spite of the advances in video frame interpolation research, benchmarking and comparing different methods remain as a nontrivial issue. Furthermore, current widely used benchmark datasets such as Middlebury [1], UCF101 [44], THUMOS 2015 [11], Vimeo-90K [52] are limited in quantity or strength of motion. While SlowFlow [18], their high frame rate version of Sintel [4] and HD dataset [3] are in higher quality, they are not much popular for frame interpolation benchmark, yet.

In this paper, we report the AIM 2019 Challenge on Video Temporal Super-Resolution with a focus on the submitted methods and the benchmark results. We provide a new large-scale dataset, REDS_VTSR to train and evaluate the video temporal super-resolution methods in a unified environment. The proposed REDS_VTSR dataset consists of high-quality dynamic scenes with large motions at different frame rates. The dataset is derived from the superset of videos that are used to create REDS [29] dataset for video deblurring [31] and super-resolution [32]. There are 30 sequences each consisting of 181 frames at 60 fps where the corresponding 30 fps and 15 fps counterparts are subsampled from. The challenge goal is to recover 60 fps video from low-frame-rate 15 fps input sequence.In the following sections, we describe the related works and introduce the AIM 2019 Video Temporal Super-Resolution challenge (VTSR). We also present and discuss the challenge results with the proposed methods.

## 2. Related Works

Classical image interpolation techniques mostly include motion modeling via optical flow [1] or phase shift [8, 50]. Modern video frame interpolation methods also employ motion estimation stemming from them, developing more sophisticated motion modeling and warping methods. Meanwhile, many video datasets are adopted for training and evaluating those proposed methods.

### Frame interpolation

Long *et al*. [26] proposed MIND, an early CNN based approach that learns to estimate the interpolated frame without explicit motion modeling. However, their primary goal was to infer the correspondence between input frames by inverting the interpolation network, and the interpolation accuracy was not evaluated.

Meyer *et al*. [28] proposed a phase-based frame interpolation by using the phase difference between input frames without using optical flow. As no direct pixel correspondence is calculated, the method is robust against the illumination change and blur. Later, Meyer *et al*. [27] proposed PhaseNet that predicts amplitude and phase decomposition of the intermediate frame. The CNN is trained in an end-to-end manner via image loss and phase loss.

There were other approaches that tried to unify motion estimation and frame synthesis. Niklaus *et al*. [34] extracted $41 \times 41$ spatially adaptive convolutional kernel from a CNN model. The kernel is convolved with the input to synthesize the intermediate frame. Niklaus *et al*. [35] factorized the spatial kernel with 1D kernels to reduce memory footprint.

On the other hand, the concept of optical flow was often embedded in neural network architectures. Liu *et al*. [25] implements a voxel flow layer where the spatial component is the optical flow at the corresponding time. By assuming the flow to be locally linear and temporally symmetric, the following volume sampling layer predicts the target frame. Niklaus *et al*. [33] relaxes the assumptions by using bidirectional flow [48] to warp the input frames and features.

However, optical flow estimation could bring errors in case of occlusion. Xue *et al*. [52] implicitly handles occlusion by learning task-specific flow followed by spatial transformer [17] for the area where warping from the optical flow may fail. In contrast, Jiang *et al*. [20] introduces a soft visibility map for occlusion reasoning to calculate the contribution of the corresponding inputs to the target time. Also, the flow interpolation network approximates the flow from the arbitrary intermediate time frame to inputs. Thus, the proposed SuperSloMo method could directly generate the intermediate frame at any intermediate moment without recursion in contrast to other methods that estimate the middle frame only. Hu et al. [13] also proposed an anytime estimation model by unifying frame interpolation and

extrapolation framework by applying transitive consistency loss.

Bao et al. [2] combines the kernel and optical flow based approaches. Their adaptive warping layer synthesizes a new pixel value by applying a local convolutional kernel where the position of the kernel window is determined by optical flow. During the warping process, relative monocular depth is estimated so that closer objects contribute more during flow projection. Moreover, Peleg et al. [37] focuses on real-time processing and a large receptive field. They extract a low-resolution feature from multi-scale architecture to acquire vertical and horizontal motion vector field. The model is trained with kernel loss [35], trilinear interpolation loss [25], etc. As their model operates on lower resolution inputs, the receptive field is large, capable of handling motion size up to $192 \times 192$.

The AIM 2019 video temporal super-resolution challenge participants mostly adopt the designs of the previous methods or modify the network architecture to improve performance or reduce computational complexity.

### Frame Interpolation Datasets

**Middlebury** [1] is one of the most popular benchmarks for frame interpolation, which was originally designed for optical flow evaluation. While it is widely used, it has only 8 scenes, limited in quantity. The metric is IE (interpolation error) which can be directly converted to PSNR.

**UCF101** [44] is a dataset designed for action classification. It contains various dynamics such as human-object interaction, sports, playing instruments. The training set was used for Deep Voxel Flow [25] while the test set is more often employed [25, 20, 52, 13, 2]. The resolution is $320 \times 240$.

**THUMOS 2015** [16] is also a dataset for action recognition whose videos are collected from YouTube. The test set contains 5613 videos of 101 action classes that are compatible with UCF101. It is used for evaluation together with the UCF101 dataset in [25, 13, 52].

**Vimeo-90k** [52] consist of $448 \times 256$ frame triplets extracted from video clips collected from Vimeo. From the originally collected pool of videos, static scenes and scenes with large illumination changes are removed. Also, nonlinear motions are also not included. However, this makes the frame interpolation task to be easier than the real-world scenario. It is used for training and testing TOF [52], DAIN [2]. IM-Net [37] uses super-resolved Vimeo sequences to evaluate the interpolation performance for large motion.

**SlowFlow** [18] are high quality dataset for optical flow. The videos are recorded with a high-speed camera at $2560 \times 1440$ resolution and $> 200$ fps. The scenes are diverse and contain varying levels of realistic blur. The **MPI Sintel** [4] are re-rendered at 1008 fps and $2048 \times 872$ resolution to check optical flow methods. They are adopted for evaluating the slow-motion algorithm [20].

| Name | Resolution | fps | # Sequences / # Frames | Note |
|---|---|---|---|---|
| Flickr [34] | $150 \times 150$ | - | 250,000 / 750,000* | *The number of training |
| YouTube [35] | $150 \times 150$ | - | 250,000 / 750,000* | patches are reported |
| YouTube [33] | $300 \times 300$ | - | 50,000 / 150,000* | instead of the number of |
| YouTube [37] | $512 \times 512$ | - | 40,000 / 120,000* | original frames. |
| YouTube + Adobe 240fps [47] | - | - | 1,132 / 376K | Used at SuperSloMo [20] |
| UCF101 [44] | $320 \times 240$ | 25 | 13,320 / - | - |
| Vimeo-90k [52] | $448 \times 256$ | - | 73,171 / 219,513 | Vimeo-90k triplet dataset |
| KITTI raw [10] (downsampled) | $384 \times 128$ | 10 | 56 / 16,951 | - |
| GoPro [30] | $1280 \times 720$ | - | 33 / 3,214 | - |
| DAVIS 2016 [5] | $720 \times 480$ | - | 50 / 3,455 | Used at PhaseNet [27] |
| DAVIS 2017 [39] | $720 \times 480$ | - | 150 / 10,459 | |
| Evaluation (test) only | | | | |
| Middlebury [1] | - | 60 | 8 / 58 | - |
| THUMOS 2015 [11, 16] | - | - | 5,613 / - | - |
| SlowFlow [18] | $1280 \times 1024$ | >200 | 46 / - | - |
| Sintel† [4, 18] | $2048 \times 872$ | 1008 | 19 / - | Synthetic data |
| HD dataset [3] | 1080p, 720p, 544p | - | 11 / - | Large motion |
| AIM 2019 VTSR challenge | | | | |
| REDS_VTSR | $1280 \times 720$ | 60 | 300 / 54,300 | Dynamic scenes |

Table 1: Frame interpolation dataset statistics. Sintel† [4] sequences are rendered in high-frame-rate [18].

**DAVIS** datasets [38, 39] are benchmark datasets for video object segmentation with a natural level of motion blur, appearance changes, camera shake, etc. They are used to train PhaseNet [27].

**KITTI raw** [10] is a dataset for autonomous driving captured on a car, recorded in 5 categories: Road, City, Residential, Campus and Person. There are 56 sequences where Long et al. [26] downsampled to $384 \times 128$. They additionaly train on downsampled Sintel [4] at $256 \times 128$.

**HD** dataset [3] are 7 video clips collected from Xiph website each consisting of 50 frames and 4 clips from Sintel [4]. The evaluation was done in [3, 2].

**GOPRO** dataset [30] is a dataset captured with a high-speed camera at 240 fps for dynamic scene deblurring. The captured frames are averaged to synthesize blurry images with reference sharp frames. MSFSN [13] is trained with it.

**Adobe 240fps** [47] dataset is also a deblurring dataset recorded with high-speed cameras. SuperSloMo [20] used it jointly with YouTube videos for training.

**Others** Several methods collected videos from the web to train their own methods. Different sets of YouTube [35, 33, 20, 37] and Flickr [34] videos are collected for each method.

As every method uses different training and test datasets, it is difficult to perform a fair comparison between them. The dataset statistics are summarized in Table 1. Also, many of the datasets lack quality in terms of resolution, diversity of motion, or quantity. In contrast, our proposed REDS_VTSR dataset contains various dynamic motion of objects and a camera. Further, we evaluate the challenge participants' methods in a unified environment.

## 3. AIM 2019 VTSR Challenge

**Challenge Goal** The AIM 2019 challenge on video temporal super-resolution is the first challenge of its kind. The purpose of the VTSR challenge is to gauge the state-of-the-art in video frame interpolation and facilitate the comparison of different solutions with a single large-scale dataset, REDS_VTSR. We show examples of the previous video datasets and the REDS_VTSR dataset in Figure 1.

The challenge objective is to recover higher-frame-rate video sequences from low-frame-rate input sequences. Given the input 15 fps videos, there were two-staged target frame rates of 30 fps and 60 fps. During the online submission period, participants submitted part of the 30 fps recovery results (1/4) due to the limitation of the evaluation server. At the final email submission period, Full 60 fps results were submitted. The results in Table 2 show the restoration performance on the full test dataset at the two frame rates.

We provide the REDS_VTSR dataset for the participants to train, validate, and test the performance. The training and validation sets are derived from the same videos as REDS [29] dataset while the test set is from a new set of videos. Each of the training, validation, and test set contains 240, 30, 30 sequences, and each sequence has 181 frames at 60 fps. The original videos are captured at 120 fps, and the 15, 30, 60 fps version of them are temporally subsampled from original videos.

To suppress the artifacts that come from high-speed recording and to avoid redundancies from REDS, we follow a similar process as [29]. We downsample the $1920 \times 1080$
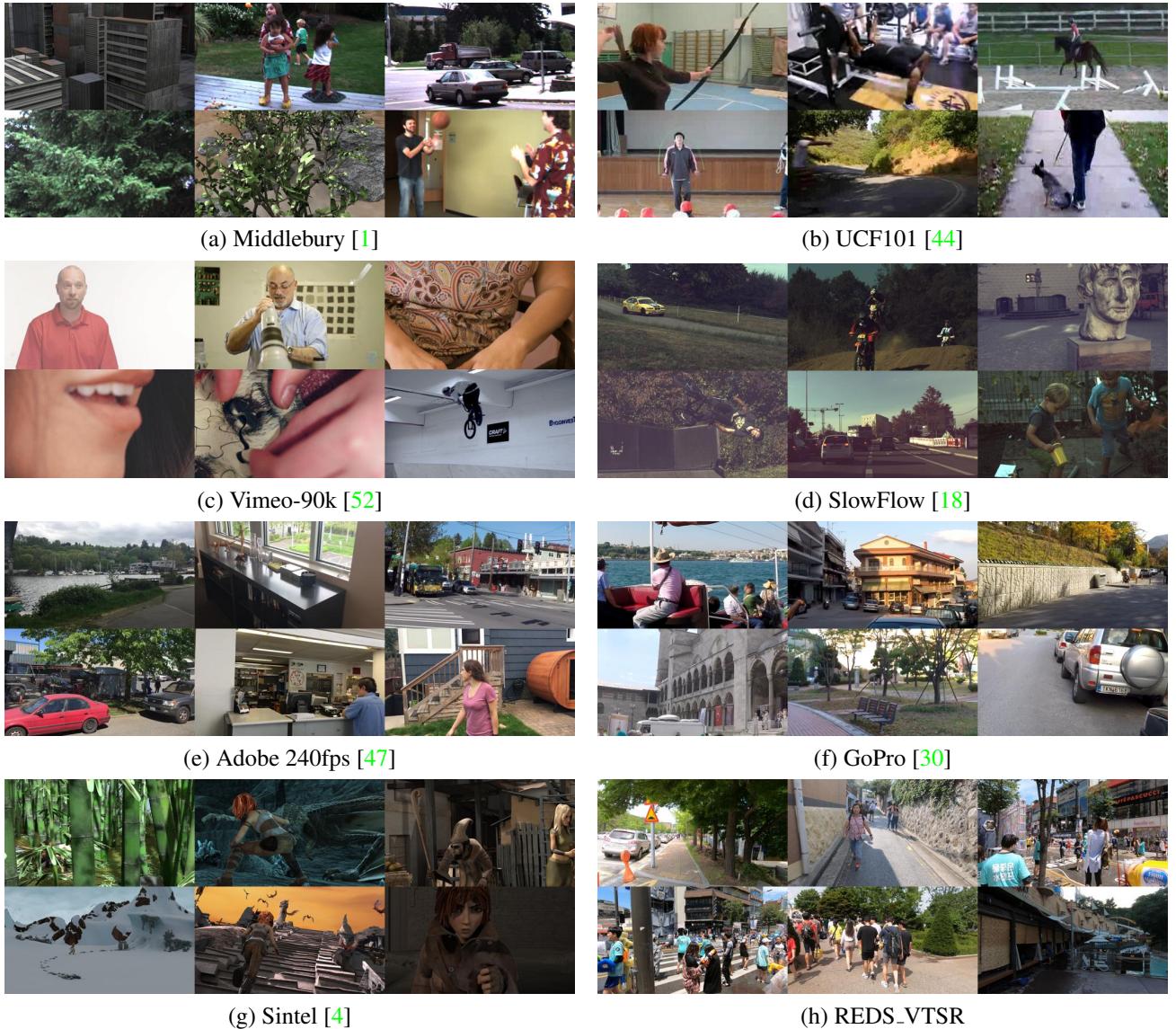
(a) Middlebury [1]

(b) UCF101 [44]

(c) Vimeo-90k [52]

(d) SlowFlow [18]

(e) Adobe 240fps [47]

(f) GoPro [30]

(g) Sintel [4]

(h) REDS_VTSR

Figure 1: Example frames from selected datasets. All samples are cropped to $16:9$ ratio for better visualization.

frames at $3/4$ ratio and center-crop to make it in standard HD resolution $1280 \times 720$. As the scene scale is relatively larger than REDS, it facilitates a relatively challenging configuration for frame interpolation with a larger strength of motion. Each frame is saved without compression.

## 4. Challenge Results

The challenge had 62 registered participants and 8 teams competed in the final test phase. The teams submitted their final 60 fps results, source code, trained models, and factsheets. All the results were reproducible from the submitted source code. We evaluated the solutions in terms of PSNR and SSIM. The competition results are summarized in Table 2 and the implementation specifications are in Table 3.

We visualize the relative performance in Figure 2.

Most of the participating teams used optical flow in their deep CNN architecture. The challenge winner, SenseSloMo [43] modeled nonlinear motion between frames with the quadratic model. They additionally use perceptual loss [21] to improve visual quality. On the other hand, NoFlow avoided using optical flow and chose to use channel attention. Interestingly, the ZSFI team uses a zero-shot approach that does not require training samples. The results of high-ranking teams are compared in Figure 3.

## 5. Challenge Methods and Teams

We describe the submitted solution details in this section. PyTorch has served as a generally favorable frame-
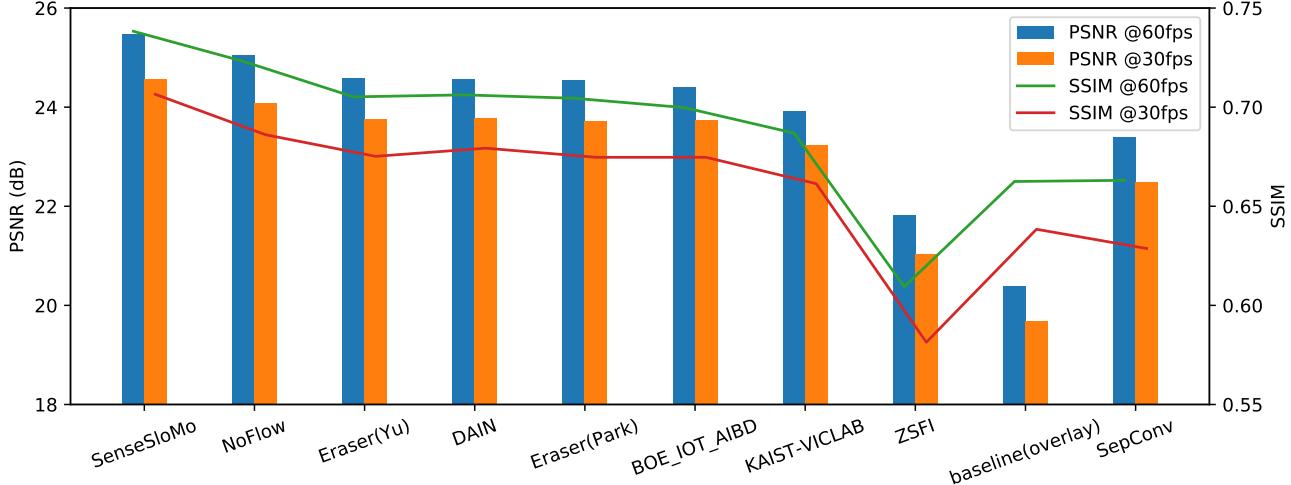
Figure 2: Test PSNR and SSIM of the solutions in 60 and 30 fps restoration.

| Team | Author | 15fps → 60fps | | 15fps → 30fps | |
|------|--------|------|------|------|------|
| | | PSNR | SSIM | PSNR | SSIM |
| SenseSloMo | Siyao | 25.47 | 0.7383 | 24.56 | 0.7065 |
| NoFlow | myungsub | 25.05 | 0.7231 | 24.08 | 0.6861 |
| Eraser | Songsaris | 24.58 | 0.7052 | 23.77 | 0.6752 |
| DAIN | wangshen233 | 24.56 | 0.7062 | 23.78 | 0.6793 |
| Eraser | BumjunPark | 24.55 | 0.7045 | 23.71 | 0.6747 |
| BOE_IOT_AIBD_IMP | BOE_IOT_AIBD_IMP | 24.41 | 0.6998 | 23.74 | 0.6747 |
| KAIST-VICLAB | WSPark | 23.93 | 0.6869 | 23.23 | 0.6614 |
| ZSFI | gpkoko | 21.83 | 0.6094 | 21.03 | 0.5814 |
| SepConv [35] | - | 23.40 | 0.6631 | 22.48 | 0.6287 |
| *baseline* (overlay) | - | 20.39 | 0.6625 | 19.68 | 0.6384 |

Table 2: AIM 2019 Video Temporal Super-Resolution Challenge results on the REDS_VTSR test data. Teams are sorted by ranking in terms of PSNR for 60 fps restoration.

| Team | Runtime (s) | Platform | GPU (at runtime) | Ensemble / Fusion (at runtime) |
|------|-------------|----------|------------------|--------------------------------|
| SenseSloMo | 1.00 | PyTorch | GTX 1060 | flip (x4) |
| NoFlow | 0.30 | PyTorch | TITAN Xp | - |
| Eraser (Yu) | 0.46 | PyTorch | TITAN V | rotation / flip (x8) |
| DAIN | 0.82 | PyTorch / ATen | RTX 2080 Ti | - |
| Eraser (Park) | 2.40 | PyTorch | RTX 2080 Ti | rotation / flip (x8) |
| BOE_IOT_AIBD_IMP | 0.26 | PyTorch / C++ | Tesla V100 | - |
| KAIST-VICLAB | 0.25 | PyTorch | TITAN Xp | - |
| ZSFI | 0.50 | TensorFlow | GTX 1080 Ti / RTX 2080 | - |

Table 3: Reported runtime per frame on REDS_VTSR test data (60fps) and details from the factsheets. We note that the ZSFI team's method requires 1-3 hour sequence-specific training.

work while some methods have leveraged customized C++ or ATen kernels for flow prediction. Several methods are reporting that run-time self-ensemble strategy [49] can enhance the model performance at the cost of running speed.

## 5.1. SenseSloMo

The SenseSloMo team takes a quadratic frame interpolation method [51, 43] which can reflect the acceleration of motions among video frames as in Fig-

| (a) GT | (b) SenseSloMo | (c) NoFlow | (d) Eraser (Yu) | (e) DAIN |

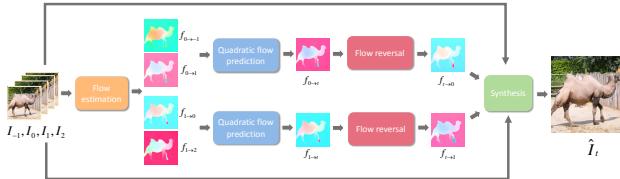Figure 3: Temporal super-resolution results (REDS_VTSR '002/00000006')



Figure 4: SenseSloMo: the proposed quadratic frame interpolation method.

ure 4. Specifically, the proposed method uses four input frames $I_{-1}, I_0, I_1, I_2$ to generate an arbitrary inter frame $I_t$ $(0 < t < 1)$. Generally, the proposed model can be divided into four parts. First, the network estimates optical flows $F_{0\to-1}, F_{0\to1}, F_{1\to0}$ and $F_{1\to2}$. Second, $F_{0\to t}$ and $F_{1\to t}$ are calculated using a quadratic flow formula. Third, the SenseSloMo team reverses those flows to $F_{t\to0}$ and $F_{t\to1}$. Finally, the interpolated frame $I_t$ is synthesized.

## 5.2. NoFlow

The NoFlow team is motivated by the potential drawbacks in using optical flow and proposes a novel framework of video frame interpolation that replaces optical flow with a simple convolutional network as in Figure 5. Encoding images to lower spatial resolution gradually distributes the motion-related information into multiple channels to con-

struct a transformed feature map. This feature representation is then combined with channel attention to capture the variations between the two input frames and synthesizes high-quality intermediate video frames without explicit motion estimation.
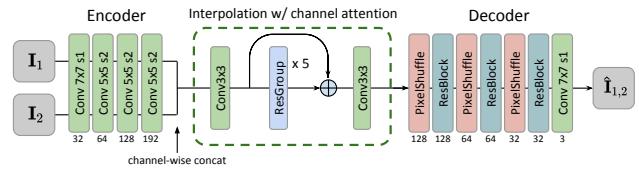


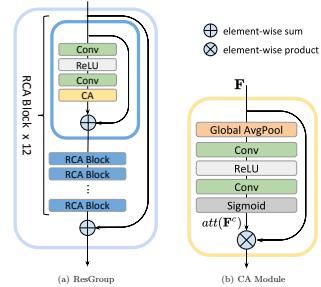Figure 5: NoFlow: overall network architecture.



Figure 6: ResGroup and CA modules in NoFlow architecture.

Using ResGroups and CA modules as in Figure 6, the proposed approach is capable of handling large motion and heavy occlusion effectively and outperforms the prior state-of-the-art methods.

## 5.3. Eraser

The Eraser team has submitted two solutions(Yu, Park) for video frame interpolation.

Eraser (Yu) predicts 3 intermediate frames for each frame pair. The frame index representation used in this approach is shown in Figure 7. The proposed method uses two separate models: a middle network to estimate a middle frame (Frame 2) and a side network to estimate side frames (Frame 1 and 3). Each model consists of optical flow estimation and enhancement as Figure 8. For initial optical flow estimation, a pre-trained PWC-Net [48] is used. Then a modified DIDN [54] is introduced as an enhancement network. In the proposed approach, at least one flow estimator exists for each frame position. For example, one estimator is used for position 1, two estimators for position 2, and one estimator for position 3. Because position 2 is in the middle, two estimators (forward direction and backward direction) can be allocated. Frames at positions 1 and 3 are predicted using only close frames (the previous
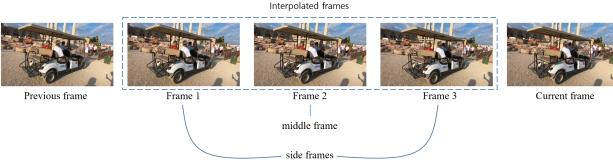


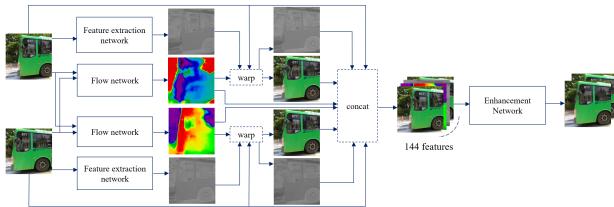Figure 7: Eraser (Yu): an example of frame index representation.



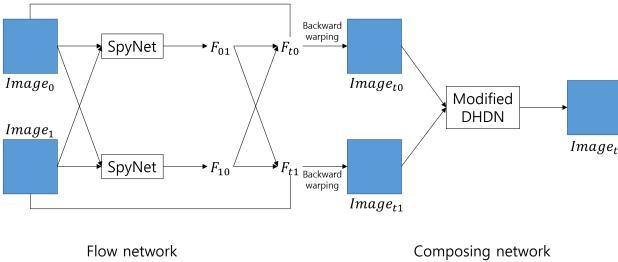Figure 8: Eraser (Yu): overall pipeline of the method.



Figure 9: Eraser (Park): the proposed DVTSR architecture.

frame is used for predicting position 1 and current frame for position 3). As an input of the enhancement network, all of the extracted features from ResNet-151 [12], warped frames from PWCNet, original frames, and flow information is used. The pre-trained flow estimator is fine-tuned for a specific frame position and is named as a network using position-specific flow (PoSNet). All networks (PWC-Net + modified DIDN) are trained together in an end-to-end manner. In other words, as in ToFlow [52], the pre-trained PWC-Net is also fine-tuned for the challenge dataset.

Eraser (Park) is inspired by TOFlow [52] and Super-SloMo [20]. Figure 9 shows the architecture of the DVTSR. The model consists of two networks: flow network and composing network. For the flow network, the Eraser team (Park) adopts the architecture of SpyNet [40]. For the composing network, the architecture of DHDN [36] is used with modification. The number of initial feature maps is 30, and one DCR block is used per each step.

## 5.4. DAIN

The DAIN team proposes a video frame interpolation method, which explicitly detects the occlusion by exploring the depth information. First, the optical flows and depth maps are extracted from video frames as features. Then, a depth-aware flow projection layer is used to synthesize intermediate flows that preferably sample closer objects than farther ones. An adaptive warping layer takes the projected flows along with depth map, encoded contextual features, and interpolation kernel features to produce warped depth maps, warped frames, and warped contextual features together. At last, the warped features are concatenated and a frame synthesis net output the resulting frame. The overall pipeline is visualized in Figure 10.
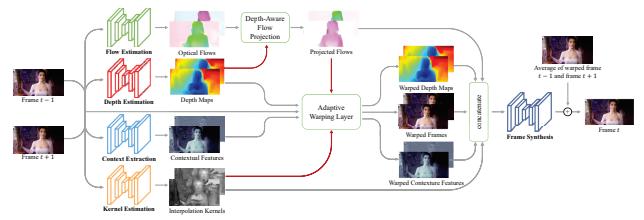


Figure 10: DAIN: overall pipeline of the method.

## 5.5. BOE_IOT_AIBD_IMP

The BOE_IOT_AIBD_IMP team has proposed the frame interpolation solution based on DAIN [2]. The team has improved the model by using IRR-PWC [15] for calculating the optical flows, PacJointUpsample [46] for joint optical flow upsampling, VNL [53] for predicting the depth maps and pixel-adaptive convolutions [46] in frame synthesis module for generating the final results. Based on

CyclicGen [24], a new two-stage cyclic generation process is proposed for training the 4x interpolation model as in Figure 11. The method is using pre-trained IRR-PWC [15], VNL [53], PacJointUpsample [46], and some sub-module/layers of DAIN [2], and was trained/finetuned on the REDS VTSR dataset. The entire network architecture is visualized in Figure 12.
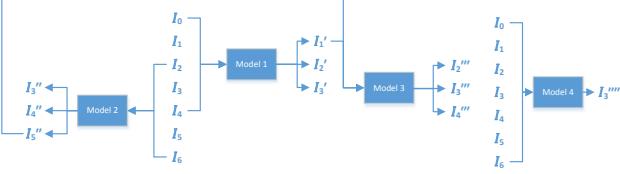


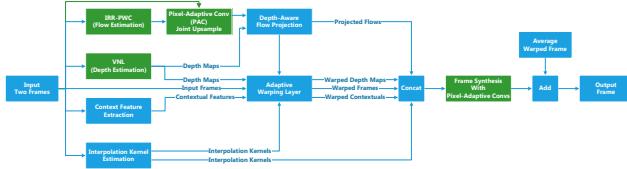Figure 11: BOE_IOT_AIBD_IMP: the proposed cyclic training process.



Figure 12: Overall network architecture of the BOE_IOT_AIBD_IMP team.

## 5.6. KAIST-VICLAB

The KAIST-VICLAB team proposes a method using SuperSloMo [20], which consists of two U-Nets as the basic structure. The original SuperSloMo structure estimates the motion vector from a single scale and performs frame interpolation. The proposed method is mainly composed of four sub-networks, as shown in Figure 13. The first sub-net (Motion Estimation) estimates the motion vectors from multi-scale (Scale 0, Scale 1, Scale 2) inputs. Then, the second sub-net (Motion Synthesis) synthesizes the motion vectors in multi-scale (Scale 0, Scale 1, Scale 2) by MVNet. The third sub-net (Motion Compensation) performs the frame interpolation for each scale using the estimated multi-scale motion vectors. Finally, the intermediate frame is obtained through the last sub-network (Image synthesis) which finally synthesizes the obtained frames for each scale. These four sub-networks of our proposed multi-scale structure is trained in an end-to-end manner. At the second sub-network (Motion Synthesis), only the residuals of motion vectors for three scales (0, 1, and 2) are calculated by our networks. Before the last sub-network (Image Synthesis), the reconstructed image for each scale is resized to the original image size by bilinear interpolation. In Scale 2, our proposed network can process the feature maps with large receptive fields, and in Scale 0, the proposed network can process textured areas efficiently.
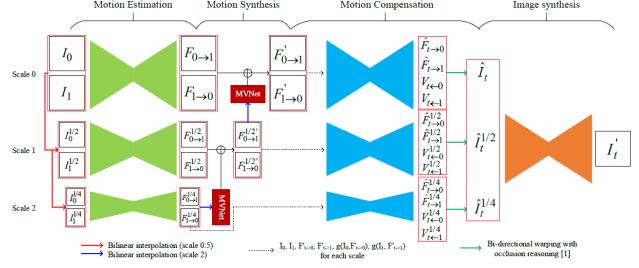


Figure 13: KAIST-VICLAB: overall pipeline of the method.

## 5.7. ZSFI

The ZSFI team proposes a novel zero-shot frame interpolation, i.e., it reconstructs and interpolates a video without any prior training. This approach is based on the concepts of internal statistics and self-similarity in the space-time volume of a video described by Shahar et al. [41]. The work has demonstrated that space-time patches recur over large distances both in space and in time and most importantly, across scales, that are in different areas of the frame and between frames. Using the concept of recurring patches, a zero-shot super-resolution algorithm [42] has proposed recently. Here, the ZSFI team extends this concept to video interpolation from a single example, the video itself as in Figure 14. From the test video, the ZSFI team creates a temporally down-sampled version of the video that allows us to leverage the internal similarity across scales and learn the interpolation-based solely on the test video. The proposed approach is especially effective when the target video is significantly different from the training data such as medical or synthetic data.
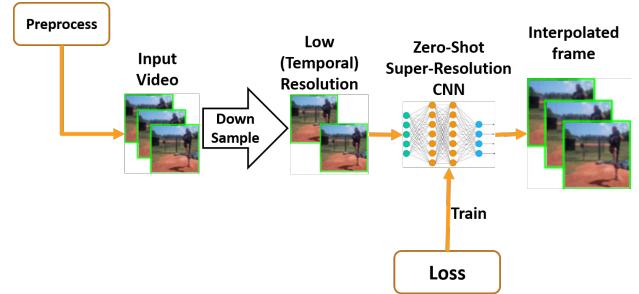


Figure 14: ZSFI: zero-shot training process.

## Acknowledgments

# A. Teams and affiliations

## AIM 2019 team

***Title:*** AIM 2019 Challenge on Video Temporal Super-Resolution
***Members:*** *Seungjun Nah[1] (seungjun.nah@gmail.com)*, Sanghyun Son[1], Radu Timofte[2], Kyoung Mu Lee[1]
***Affiliations:***
[1] Department of ECE, ASRI, Seoul National University (SNU), Korea
[2] Computer Vision Lab, ETH Zurich, Switzerland

## SenseSloMo

***Title:*** Quadratic Video Interpolation
***Members:*** *Li Siyao (lisiyao1@sensetime.com)*, Ze Pan, Xiangyu Xu, Wenxiu Sun
***Affiliations:***
SenseTime Research, China

## NoFlow

***Title:*** Channel Attention is All You Need for Video Temporal Super-Resolution
***Members:*** *Myungsub Choi[1] (cms6539@gmail.com)*, Heewon Kim[1], Bohyung Han[1], Ning Xu[2], Kyoung Mu Lee[1]
***Affiliations:***
[1] Department of ECE, ASRI, Seoul National University (SNU), Korea
[2] Amazon, USA

## Eraser

***Title:*** Video Frame Interpolation Using Position-Specific Flow
***Members:*** *Bumjun Park (kkbbbj@gmail.com)*, Songhyun Yu, Sangmin Kim, Jechang Jeong
***Affiliations:***
Department of ECE, Hanyang University, Korea

## DAIN

***Title:*** Depth-Aware Video Frame Interpolation
***Members:*** *Wang Shen (shenwang@sjtu.edu.cn)*, Wenbo Bao, Guangtao Zhai, Li Chen, Zhiyong Gao
***Affiliations:***
Department of EE, Shanghai Jiao Tong University, China

## BOE_IOT_AIBD_IMP

***Title:*** Pixel-Adaptive Joint Depth-aware Cyclic Network for Video Temporal Super-resolution
***Members:*** *Guannan Chen (578489493@qq.com)*, Yunhua Lu, Ran Duan, Tong Liu, Lijie Zhang
***Affiliations:***
BOE Technology Group Co., Ltd., China

## KAIST-VICLAB

***Title:*** Multi-scale Motion Estimation and Motion Compensation
***Members:*** *Woonsung Park (pys5309@kaist.ac.kr)*, Munchurl Kim
***Affiliations:***
Korea Advanced Institute of Science and Technology (KAIST), Korea

## ZSFI

***Title:*** Zero Shot Frame Interpolation
***Members:*** *George Pisha (pisha@campus.technion.ac.il)*, Eyal Naor, Lior Aloni
***Affiliations:***
Technion Israel Institute of Technology, Israel

## References

[1] S. Baker, D. Scharstein, J. P. Lewis, S. Roth, M. J. Black, and R. Szeliski. A database and evaluation methodology for optical flow. *International Journal of Computer Vision (IJCV)*, 92(1):1–31, Mar. 2011. 1, 2, 3, 4

[2] W. Bao, W.-S. Lai, C. Ma, X. Zhang, Z. Gao, and M.-H. Yang. Depth-aware video frame interpolation. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Jun. 2019. 1, 2, 3, 7, 8

[3] W. Bao, W.-S. Lai, X. Zhang, Z. Gao, and M.-H. Yang. MEMC-Net: Motion estimation and motion compensation driven neural network for video interpolation and enhancement. *arXiv preprint arXiv:1810.08768*, 2018. 1, 3

[4] D. J. Butler, J. Wulff, G. B. Stanley, and M. J. Black. A naturalistic open source movie for optical flow evaluation. In *The European Conference on Computer Vision (ECCV)*, Oct. 2012. 1, 2, 3, 4

[5] S. Caelles, K.-K. Maninis, J. Pont-Tuset, L. Leal-Taixé, D. Cremers, and L. Van Gool. One-shot video object segmentation. In *The IEEE conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. 3

[6] R. Castagno, P. Haavisto, and G. Ramponi. A method for motion adaptive frame rate up-conversion. *IEEE Transactions on circuits and Systems for Video Technology*, 6(5):436–446, 1996. 1

[7] B.-D. Choi, J.-W. Han, C.-S. Kim, and S.-J. Ko. Motion-compensated frame interpolation using bilateral motion estimation and adaptive overlapped block motion compensation. *IEEE Transactions on Circuits and Systems for Video Technology*, 17(4):407–416, 2007. 1

[8] P. Didyk, P. Sitthi-Amorn, W. Freeman, F. Durand, and W. Matusik. Joint view expansion and filtering for automul-

tiscopic 3d displays. *ACM Transactions on Graphics (TOG)*, 32(6):221, 2013. 2

[9] J. Flynn, I. Neulander, J. Philbin, and N. Snavely. Deep-stereo: Learning to predict new views from the world's imagery. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Jun. 2016. 1

[10] A. Geiger, P. Lenz, C. Stiller, and R. Urtasun. Vision meets robotics: The KITTI dataset. *International Journal of Robotics Research (IJRR)*, 2013. 1, 3

[11] A. Gorban, H. Idrees, Y.-G. Jiang, A. Roshan Zamir, I. Laptev, M. Shah, and R. Sukthankar. THUMOS challenge: Action recognition with a large number of classes. http://www.thumos.info/, 2015. 1, 3

[12] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Jun. 2016. 7

[13] Z. Hu, Y. Ma, and L. Ma. Multi-scale video frame-synthesis network with transitive consistency loss. *arXiv preprint arXiv:1712.02874*, 2017. 1, 2, 3

[14] X. Huang, L. L. Rakêt, H. Van Luong, M. Nielsen, F. Lauze, and S. Forchhammer. Multi-hypothesis transform domain wyner-ziv video coding including optical flow. In *2011 IEEE 13th International Workshop on Multimedia Signal Processing*, 2011. 1

[15] J. Hur and S. Roth. Iterative residual refinement for joint optical flow and occlusion estimation. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Jun. 2019. 7, 8

[16] H. Idrees, A. R. Zamir, Y. Jiang, A. Gorban, I. Laptev, R. Sukthankar, and M. Shah. The THUMOS challenge on action recognition for videos in the wild. *Computer Vision and Image Understanding (CVIU)*, 155:1–23, 2017. 2, 3

[17] M. Jaderberg, K. Simonyan, A. Zisserman, and k. kavukcuoglu. Spatial transformer networks. In *Advances in Neural Information Processing Systems (NIPS)*. 2015. 2

[18] J. Janai, F. Gney, J. Wulff, M. Black, and A. Geiger. Slow Flow: Exploiting high-speed cameras for accurate and diverse optical flow reference data. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Jun. 2017. 1, 2, 3, 4

[19] B.-W. Jeon, G.-I. Lee, S.-H. Lee, and R.-H. Park. Coarse-to-fine frame interpolation for frame rate up-conversion using pyramid structure. *IEEE Transactions on Consumer Electronics*, 49(3):499–508, 2003. 1

[20] H. Jiang, D. Sun, V. Jampani, M.-H. Yang, E. Learned-Miller, and J. Kautz. Super SloMo: High quality estimation of multiple intermediate frames for video interpolation. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Jun. 2018. 1, 2, 3, 7, 8

[21] J. Johnson, A. Alahi, and L. Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *The European Conference on Computer Vision (ECCV)*, pages 694–711. Springer, 2016. 4

[22] S.-J. Kang, K.-R. Cho, and Y. H. Kim. Motion compensated frame rate up-conversion using extended bilateral motion estimation. *IEEE Transactions on Consumer Electronics*, 53(4):1759–1767, 2007. 1

[23] S.-H. Lee, O. Kwon, and R.-H. Park. Weighted-adaptive motion-compensated frame rate up-conversion. *IEEE Transactions on Consumer Electronics*, 49(3):485–492, 2003. 1

[24] Y.-L. Liu, Y.-T. Liao, Y.-Y. Lin, and Y.-Y. Chuang. Deep video frame interpolation using cyclic frame generation. In *Proceedings of the 33rd Conference on Artificial Intelligence (AAAI)*, 2019. 8

[25] Z. Liu, R. A. Yeh, X. Tang, Y. Liu, and A. Agarwala. Video frame synthesis using deep voxel flow. In *The IEEE International Conference on Computer Vision (ICCV)*, Oct. 2017. 1, 2

[26] G. Long, L. Kneip, J. M. Alvarez, H. Li, X. Zhang, and Q. Yu. Learning image matching by simply watching video. In *The European Conference on Computer Vision (ECCV)*, 2016. 1, 2, 3

[27] S. Meyer, A. Djelouah, B. McWilliams, A. Sorkine-Hornung, M. Gross, and C. Schroers. Phasenet for video frame interpolation. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Jun. 2018. 1, 2, 3

[28] S. Meyer, O. Wang, H. Zimmer, M. Grosse, and A. Sorkine-Hornung. Phase-Based frame interpolation for video. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Jun. 2015. 2

[29] S. Nah, S. Baik, S. Hong, G. Moon, S. Son, R. Timofte, and K. M. Lee. NTIRE 2019 challenges on video deblurring and super-resolution: Dataset and study. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, Jun. 2019. 1, 3

[30] S. Nah, T. H. Kim, and K. M. Lee. Deep multi-scale convolutional neural network for dynamic scene deblurring. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Jul. 2017. 1, 3, 4

[31] S. Nah, R. Timofte, S. Baik, S. Hong, G. Moon, S. Son, K. M. Lee, X. Wang, K. C. Chan, K. Yu, C. Dong, C. C. Loy, Y. Fan, J. Yu, D. Liu, T. S. Huang, H. Sim, M. Kim, D. Park, J. Kim, S. Y. Chun, M. Haris, G. Shakhnarovich, N. Ukita, S. W. Zamir, A. Arora, S. Khan, F. S. Khan, L. Shao, R. K. Gupta, V. Chudasama, H. Patel, K. Upla, H. Fan, G. Li, Y. Zhang, X. Li, W. Zhang, Q. He, K. Purohit, A. N. Rajagopalan, J. Kim, M. Tofighi, T. Guo, and V. Monga. NTIRE 2019 challenge on video deblurring: Methods and results. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, Jun. 2019. 1

[32] S. Nah, R. Timofte, S. Gu, S. Baik, S. Hong, G. Moon, S. Son, K. M. Lee, X. Wang, K. C. Chan, K. Yu, C. Dong, C. C. Loy, Y. Fan, J. Yu, D. Liu, T. S. Huang, X. Liu, C. Li, D. He, Y. Ding, S. Wen, F. Porikli, R. Kalarot, M. Haris, G. Shakhnarovich, N. Ukita, P. Yi, Z. Wang, K. Jiang, J. Jiang, J. Ma, H. Dong, X. Zhang, Z. Hu, K. Kim, D. U. Kang, S. Y. Chun, K. Purohit, A. N. Rajagopalan, Y. Tian, Y. Zhang, Y. Fu, C. Xu, A. M. Tekalp, M. A. Yilmaz, C. Korkmaz, M. Sharma, M. Makwana, A. Badhwar, A. P. Singh, A. Upadhyay, R. Mukhopadhyay, A. Shukla, D. Khanna, A. Mandal, S. Chaudhury, S. Miao, Y. Zhu, and X. Huo. NTIRE 2019 challenge on video super-resolution: Methods and results. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, Jun. 2019. 1

[33] S. Niklaus and F. Liu. Context-aware synthesis for video frame interpolation. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Jun. 2018. 1, 2, 3

[34] S. Niklaus, L. Mai, and F. Liu. Video frame interpolation via adaptive convolution. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Jul. 2017. 1, 2, 3

[35] S. Niklaus, L. Mai, and F. Liu. Video frame interpolation via adaptive separable convolution. In *The IEEE International Conference on Computer Vision (ICCV)*, Oct. 2017. 1, 2, 3, 5

[36] B. Park, S. Yu, and J. Jeong. Densely connected hierarchical network for image denoising. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, Jun. 2019. 7

[37] T. Peleg, P. Szekely, D. Sabo, and O. Sendik. Im-net for high resolution video frame interpolation. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Jun. 2019. 1, 2, 3

[38] F. Perazzi, J. Pont-Tuset, B. McWilliams, L. Van Gool, M. Gross, and A. Sorkine-Hornung. A benchmark dataset and evaluation methodology for video object segmentation. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Jun. 1, 2

[39] J. Pont-Tuset, F. Perazzi, S. Caelles, P. Arbeláez, A. Sorkine-Hornung, and L. Van Gool. The 2017 DAVIS challenge on video object segmentation. *arXiv preprint arXiv:1704.00675*, 2017. 1, 2, 3

[40] A. Ranjan and M. J. Black. Optical flow estimation using a spatial pyramid network. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Jul. 2017. 7

[41] O. Shahar, A. Faktor, and M. Irani. Space-Time super-resolution from a single video. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Jun. 2011. 8

[42] A. Shocher, N. Cohen, and M. Irani. zero-shot super-resolution using deep internal learning. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Jun. 2018. 8

[43] L. Siyao, X. Xu, Z. Pan, and W. Sun. Quadratic video interpolation for vtsr challenge. In *The IEEE International Conference on Computer Vision (ICCV)*. 4, 5

[44] K. Soomro, A. R. Zamir, and M. Shah. UCF101: A dataset of 101 human actions classes from videos in the wild. *CRCV-TR-12-01*, Nov. 2012. 1, 2, 3, 4

[45] T. Stich, C. Linz, G. Albuquerque, and M. Magnor. View and time interpolation in image space. In *Computer Graphics Forum*, volume 27, pages 1781–1787. Wiley Online Library, 2008. 1

[46] H. Su, V. Jampani, D. Sun, O. Gallo, E. Learned-Miller, and J. Kautz. Pixel-Adaptive convolutional neural networks. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Jun. 2019. 7, 8

[47] S. Su, M. Delbracio, J. Wang, G. Sapiro, W. Heidrich, and O. Wang. Deep video deblurring for hand-held cameras. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Jul. 2017. 1, 3, 4

[48] D. Sun, X. Yang, M.-Y. Liu, and J. Kautz. PWC-Net: CNNs for optical flow using pyramid, warping, and cost volume. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Jun. 2018. 2, 7

[49] R. Timofte, R. Rothe, and L. Van Gool. Seven ways to improve example-based single image super resolution. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Jun. 2016. 5

[50] N. Wadhwa, M. Rubinstein, F. Durand, and W. T. Freeman. Phase-based video motion processing. *ACM Transactions on Graphics (TOG)*, 32(4):80, 2013. 2

[51] X. Xu, L. Siyao, W. Sun, Q. Yin, and M.-H. Yang. Quadratic video interpolation. In *NeurIPS*, 2019. 5

[52] T. Xue, B. Chen, J. Wu, D. Wei, and W. T. Freeman. Video enhancement with task-oriented flow. *International Journal of Computer Vision (IJCV)*, 127(8):1106–1125, 2019. 1, 2, 3, 4, 7

[53] W. Yin, Y. Liu, C. Shen, and Y. Yan. Enforcing geometric constraints of virtual normal for depth prediction. In *The IEEE International Conference on Computer Vision (ICCV)*, 2019. 7, 8

[54] S. Yu, B. Park, and J. Jeong. Deep iterative down-up cnn for image denoising. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, Jun. 2019. 7