

Asymmetric Bilateral Motion Estimation for Video Frame Interpolation

Junheum Park
Korea University
jhpark@mcl.korea.ac.kr

Chul Lee
Dongguk University
chullee@dongguk.edu

Chang-Su Kim
Korea University
changasukim@korea.ac.kr

Abstract

We propose a novel video frame interpolation algorithm based on asymmetric bilateral motion estimation (ABME), which synthesizes an intermediate frame between two input frames. First, we predict symmetric bilateral motion fields to interpolate an anchor frame. Second, we estimate asymmetric bilateral motions fields from the anchor frame to the input frames. Third, we use the asymmetric fields to warp the input frames backward and reconstruct the intermediate frame. Last, to refine the intermediate frame, we develop a new synthesis network that generates a set of dynamic filters and a residual frame using local and global information. Experimental results show that the proposed algorithm achieves excellent performance on various datasets. The source codes and pretrained models are available at <https://github.com/JunHeum/ABME>.

1. Introduction

Video frame interpolation is a low-level vision task to increase the frame rate of a video sequence by interpolating intermediate frames between successive input frames. It is widely used in applications, including video enhancement [42], video compression [25], slow-motion generation [16], and view synthesis [11, 18]. Due to its practical importance, various algorithms have been proposed to increase video frame rates [2, 3, 6–8, 13, 15, 16, 20, 22–24, 30–34].

These algorithms can be classified into three categories: kernel-based [2, 3, 6, 20, 32, 33], phase-based [27, 28], and motion-based [2, 3, 13, 16, 22, 23, 30, 31, 34]. With the recent advances in optical flow estimation [9, 14, 17, 21, 35, 39, 43], motion-based algorithms have been developed most actively. They use optical flows to predict an intermediate frame by warping two successive frames forward or backward. For example, Niklaus and Liu [30] predict bi-directional optical flows between two frames and halve them to generate intermediate frames based on forward warping. However, forward warping may cause interpolation artifacts in holes and overlapped regions [41]. To overcome the hole issue, they develop a synthesis network that

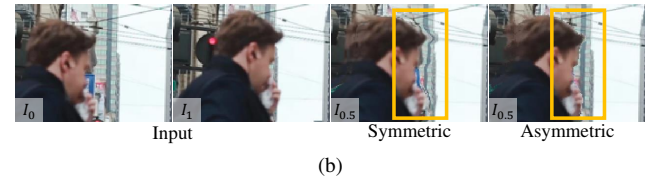
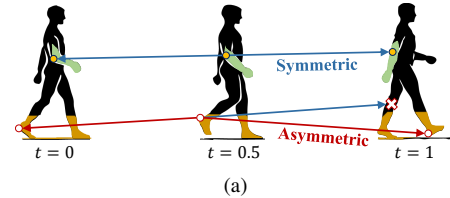


Figure 1. Symmetric vs. asymmetric bilateral motion models. In (a), the asymmetric model represents bilateral motion vectors from an intermediate frame $I_{0.5}$ to two input frames I_0 and I_1 accurately, where the symmetric one fails, by loosening the linear constraint. In (b), when $I_{0.5}$ is interpolated from I_0 and I_1 , the asymmetric model provides more faithful reconstruction with less artifacts, especially around the head, than the symmetric one does.

learns to fill in holes. The overlapping issue, however, remains, so they propose softmax-splatting [31] to combine overlapping pixel information adaptively and render the target pixel more faithfully.

On the other hand, many algorithms [2, 3, 7, 16, 22, 30, 42] are based on backward warping, which is free from the hole and overlapping issues. Backward warping needs motion vectors from intermediate frames to input frames, but intermediate frames, which should be interpolated, are unavailable at the time of motion estimation. Thus, conventional algorithms [2, 3, 16] approximate those intermediate motion vectors using optical flows between input frames. However, the approximation errors may degrade frame interpolation performance. Park *et al.* [34] employ the symmetric bilateral motion estimation to improve the accuracy of the intermediate motion, assuming that motion trajectories between input frames are linear. The linear motion constraint, however, may cause inaccurate motion estimation in regions where the constraint is invalid, such as occluded regions around motion boundaries, as illustrated in Figure 1.

In this paper, we propose a novel video frame interpola-

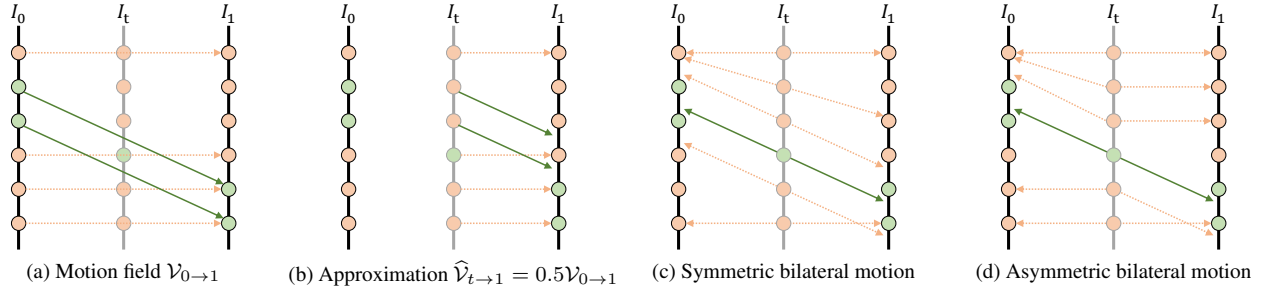


Figure 2. Illustration of various motion fields: Each column represents a frame, and a dot corresponds to a pixel in the frame. I_0 and I_1 are input frames, and I_t is an unavailable intermediate frame at time instance $t = 0.5$. Orange dots depict background pixels without any movement, while green dots depict moving objects. The motion field $\mathcal{V}_{0 \rightarrow 1}$ in (a) is halved to approximate the motion field $\hat{\mathcal{V}}_{t \rightarrow 1}$ in (b). The symmetric bilateral motion fields in (c) can be estimated to interpolate I_t based on backward warping. To improve the video frame interpolation performance, we propose the ABME algorithm, illustrated in (d).

tion algorithm based on backward warping, composed of the asymmetric bilateral motion estimation (ABME) and the frame synthesis network. In ABME, we predict symmetric bilateral motion fields and refine them by loosening the linear motion constraint. Specifically, we interpolate a temporary intermediate frame, called an anchor frame, using the symmetric fields. Then, we estimate asymmetric bilateral motion fields from the anchor frame to the two input frames, as illustrated by the red arrows in Figure 1. In the frame synthesis, the input frames are warped using the bilateral motion fields. To aggregate these warped frames, we develop a synthesis network composed of two subnetworks: FilterNet and RefineNet. FilterNet generates dynamic filters to exploit local information, while RefineNet reconstructs a residual frame using global information. Experimental results demonstrate that the proposed ABME algorithm outperforms the state-of-the-art video interpolators [2, 8, 20, 22, 33, 34, 42] meaningfully on various datasets.

2. Motion-Based Frame Warping

Let us review motion-based frame warping techniques for video frame interpolation, and introduce the notations and concepts necessary to describe how the proposed ABME is different from the conventional techniques.

Forward and backward warping: Let $\mathcal{V}_{S \rightarrow T}$ denote a pixel-wise motion field (or optical flow) from a source frame I_S to a target frame I_T . Then, the target frame can be approximated by forward warping the source frame [10],

$$\hat{I}_T = \phi_F(I_S, \mathcal{V}_{S \rightarrow T}) \quad (1)$$

where ϕ_F is the forward warping operator. On the contrary, the source frame can be approximated by backward warping the target frame [40],

$$\hat{I}_S = \phi_B(\mathcal{V}_{S \rightarrow T}, I_T) \quad (2)$$

where ϕ_B is the backward warping operator.

Given two input frames I_0 and I_1 at adjacent time instances 0 and 1, video frame interpolation aims to synthesize an intermediate frame I_t , where $0 < t < 1$. This can be achieved, using the forward warping, by

$$\hat{I}_{t,F} = (1-t) \cdot \phi_F(I_0, \mathcal{V}_{0 \rightarrow t}) + t \cdot \phi_F(I_1, \mathcal{V}_{1 \rightarrow t}). \quad (3)$$

Here, the required motion fields are often obtained by scaling the motion fields $\mathcal{V}_{0 \rightarrow 1}$ and $\mathcal{V}_{1 \rightarrow 0}$ between the input frames [30, 31], which are given by

$$\mathcal{V}_{0 \rightarrow t} = t \cdot \mathcal{V}_{0 \rightarrow 1} \quad (4)$$

$$\mathcal{V}_{1 \rightarrow t} = (1-t) \cdot \mathcal{V}_{1 \rightarrow 0}. \quad (5)$$

However, as illustrated in Figure 2(a), the scaled $\mathcal{V}_{0 \rightarrow t}$ does not pass through pixels in I_t exactly in general. Moreover, no flow vector may pass near a certain pixel, or multiple vectors may pass near the same pixel, causing hole or occlusion problems, respectively. Softmax splatting [31] alleviates these problems in the forward warping. On the other hand, a majority of video frame interpolation methods [2, 3, 16, 22, 23, 34, 36, 42], as well as the proposed algorithm, instead use the backward warping,

$$\hat{I}_{t,B} = (1-t) \cdot \phi_B(\mathcal{V}_{t \rightarrow 0}, I_0) + t \cdot \phi_B(\mathcal{V}_{t \rightarrow 1}, I_1). \quad (6)$$

However, unlike (4) and (5), it is not straightforward to obtain the motion fields $\mathcal{V}_{t \rightarrow 0}$ and $\mathcal{V}_{t \rightarrow 1}$ because the intermediate frame I_t is unavailable.

Motion approximation for backward warping: Conventional algorithms [2, 3, 16, 34, 36] approximate the motion fields $\mathcal{V}_{t \rightarrow 0}$ and $\mathcal{V}_{t \rightarrow 1}$ in (6). For example, the flow projection in [2, 3] approximates $\mathcal{V}_{t \rightarrow 0}$ and $\mathcal{V}_{t \rightarrow 1}$ by aggregating multiple flow vectors between I_0 and I_1 , which pass near each pixel in I_t . Alternatively, some algorithms simply borrow flow vectors from $\mathcal{V}_{0 \rightarrow 1}$ and $\mathcal{V}_{1 \rightarrow 0}$ to approximate $\mathcal{V}_{t \rightarrow 0}$ and $\mathcal{V}_{t \rightarrow 1}$ [34],

$$\hat{\mathcal{V}}_{t \rightarrow 0} = -t \cdot \mathcal{V}_{0 \rightarrow 1} \quad \text{or} \quad t \cdot \mathcal{V}_{1 \rightarrow 0} \quad (7)$$

$$\hat{\mathcal{V}}_{t \rightarrow 1} = (1-t) \cdot \mathcal{V}_{0 \rightarrow 1} \quad \text{or} \quad -(1-t) \cdot \mathcal{V}_{1 \rightarrow 0}. \quad (8)$$

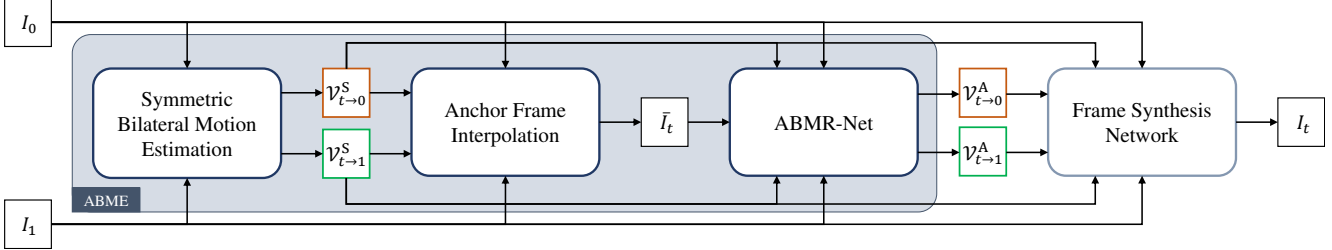


Figure 3. An overview of the proposed algorithm. ABMR-Net is detailed in Figure 5, and the frame synthesis network in Figure 4.

In [16, 36], the motion fields are approximated by combining the candidates in (7) and (8), given by

$$\widehat{\mathcal{V}}_{t \rightarrow 0} = -(1-t)t \cdot \mathcal{V}_{0 \rightarrow 1} + t^2 \cdot \mathcal{V}_{1 \rightarrow 0} \quad (9)$$

$$\widehat{\mathcal{V}}_{t \rightarrow 1} = (1-t)^2 \cdot \mathcal{V}_{0 \rightarrow 1} - t(1-t) \cdot \mathcal{V}_{1 \rightarrow 0}. \quad (10)$$

These approximations in (7)~(10) assume that neighboring pixels have similar motion vectors. However, as in Figure 2(b), the assumption is invalid around motion boundaries. In such cases, the qualities of the approximate motion fields are degraded, resulting in poorly interpolated frames.

Symmetric bilateral motion estimation: Instead of approximating the bilateral motion fields $\mathcal{V}_{t \rightarrow 0}$ and $\mathcal{V}_{t \rightarrow 1}$ using the fields $\mathcal{V}_{0 \rightarrow 1}$ and $\mathcal{V}_{1 \rightarrow 0}$ between the input frames, Park *et al.* [34] proposed the symmetric bilateral motion estimation algorithm, assuming that a motion trajectory between I_0 and I_1 is linear. Under the linear assumption, the bilateral motion fields $\mathcal{V}_{t \rightarrow 0}$ and $\mathcal{V}_{t \rightarrow 1}$ are symmetric with respect to I_t , as in Figure 2(c). More specifically,

$$\mathcal{V}_{t \rightarrow 0} = -\frac{t}{1-t} \mathcal{V}_{t \rightarrow 1}. \quad (11)$$

Therefore, roughly speaking, they obtained $\mathcal{V}_{t \rightarrow 1}$ to minimize the frame difference

$$\|\phi_B(\mathcal{V}_{t \rightarrow 0}, I_0) - \phi_B(\mathcal{V}_{t \rightarrow 1}, I_1)\| = \|\phi_B(-\frac{t}{1-t} \mathcal{V}_{t \rightarrow 1}, I_0) - \phi_B(\mathcal{V}_{t \rightarrow 1}, I_1)\|. \quad (12)$$

To this end, they developed the bilateral motion network with the bilateral cost volume.

3. Proposed Algorithm

The proposed algorithm is composed of two procedures: ABME and frame synthesis.

3.1. ABME

In Figure 2(c), bilateral motion vectors convey valid motion information between symmetrically matched pixel pairs in input frames. However, when a pixel in I_t is occluded in either I_0 or I_1 , the symmetry does not hold. Non-linear object motions due to acceleration also break the symmetry, as in Figure 1(a). To overcome these issues, we

develop the ABME technique, which refines symmetric bilateral motion vectors so that they become asymmetric and represent motion more reliably and more accurately.

Figure 3 presents an overview of the proposed algorithm. We first obtain symmetric bilateral motion fields $\mathcal{V}_{t \rightarrow 0}^S$ and $\mathcal{V}_{t \rightarrow 1}^S$ by employing the motion estimator of BMBC [34]. Using $\mathcal{V}_{t \rightarrow 0}^S$ and $\mathcal{V}_{t \rightarrow 1}^S$, we interpolate an anchor frame \bar{I}_t , which is then used as a source frame for the asymmetric bilateral motion refinement (ABMR). Finally, we obtain asymmetric bilateral motion fields $\mathcal{V}_{t \rightarrow 0}^A$ and $\mathcal{V}_{t \rightarrow 1}^A$.

Anchor frame interpolation: The motion estimation from a source frame I_t to a target frame I_0 or I_1 is challenging, since I_t is unavailable and should be synthesized in the video frame interpolation. Hence, we generate a temporary source frame \bar{I}_t , called an anchor frame, using the symmetric bilateral motion fields $\mathcal{V}_{t \rightarrow 0}^S$ and $\mathcal{V}_{t \rightarrow 1}^S$.

Based on the backward warping in (6), we may generate

$$\bar{I}_t = (1-t) \cdot \phi_B(\mathcal{V}_{t \rightarrow 0}^S, I_0) + t \cdot \phi_B(\mathcal{V}_{t \rightarrow 1}^S, I_1). \quad (13)$$

This simple blending, however, may cause errors due to occlusion, especially in boundary regions of the anchor frame in the case of camera panning. To reduce such errors, we exploit masks to reveal occluded regions, given by

$$M_{t \rightarrow 0}^S = \phi_B(\mathcal{V}_{t \rightarrow 0}^S, \mathbf{1}) \text{ and } M_{t \rightarrow 1}^S = \phi_B(\mathcal{V}_{t \rightarrow 1}^S, \mathbf{1}) \quad (14)$$

where $\mathbf{1}$ is a binary image of all ones. Note that a mask value 0 means that the corresponding pixel in I_t moves outside the frame at time instance 0 or 1. We then reconstruct the anchor frame in an occlusion-aware manner,

$$\begin{aligned} \bar{I}_t &= (1-t) \cdot (\mathbf{1} - M_{t \rightarrow 1}^S + M_{t \rightarrow 0}^S) \otimes \phi_B(\mathcal{V}_{t \rightarrow 0}^S, I_0) \\ &\quad + t \cdot (\mathbf{1} - M_{t \rightarrow 0}^S + M_{t \rightarrow 1}^S) \otimes \phi_B(\mathcal{V}_{t \rightarrow 1}^S, I_1) \end{aligned} \quad (15)$$

where \otimes is the Hadamard product.

Asymmetric bilateral motion refinement: To perform the backward warping in (6), conventional algorithms approximate bilateral motion fields $\mathcal{V}_{t \rightarrow 0}$ and $\mathcal{V}_{t \rightarrow 1}$ using the motion fields between I_0 and I_1 . In contrast, we estimate the motion field from I_t to I_0 or I_1 directly, after approximating I_t with the anchor frame \bar{I}_t .

Let us describe the asymmetric motion estimation from \bar{I}_t to I_1 . Note that the estimation from \bar{I}_t to I_0 is performed

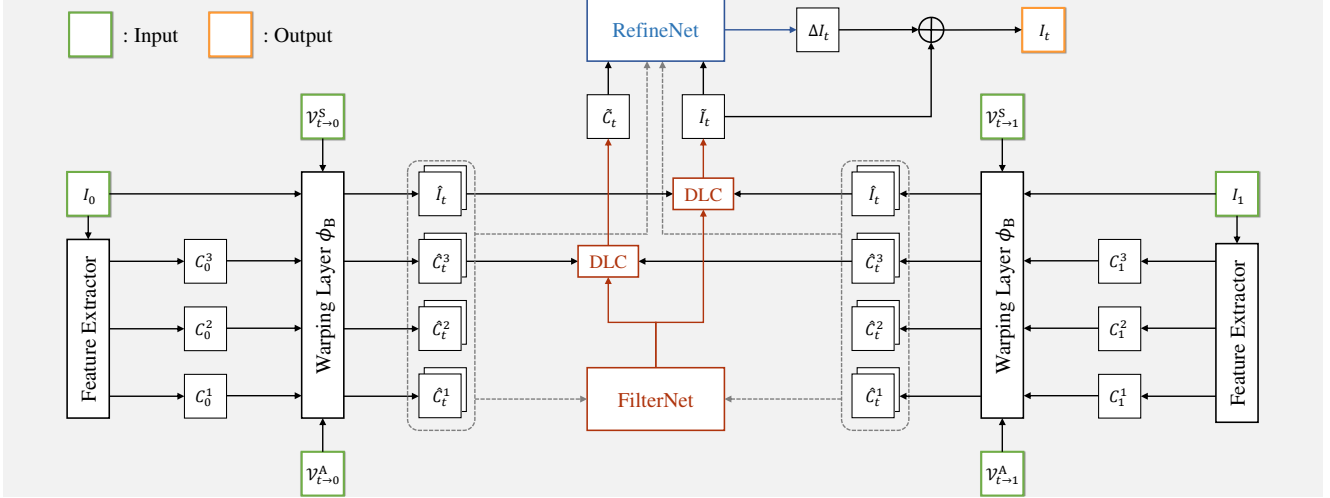


Figure 4. Architecture of the proposed frame synthesis network.

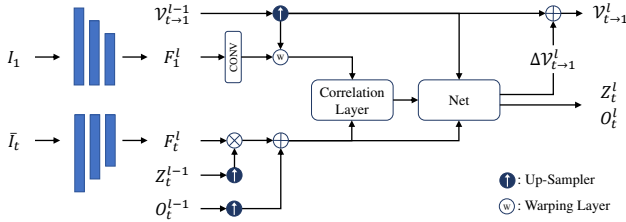


Figure 5. Architecture of ABMR-Net.

similarly but independently. We develop ABMR-Net for asymmetric bilateral motion refinement in Figure 5 to refine the motion field from the source \bar{I}_t to the target I_1 . ABMR-Net hierarchically obtains the motion field, as done in PWC-Net [39]. At level l , the motion field $\mathcal{V}_{t \rightarrow 1}^{l-1}$ at the previous level ($l - 1$) is up-sampled to warp the target feature map F_t^l . Also, we multiply the anchor feature map F_t^l with a reliability mask Z_t^{l-1} and compensate for masked-out features by adding an offset map O_t^{l-1} . Then, the warped target feature map and the compensated anchor feature map are input to the correlation layer to compute matching costs. The cost volume is used to generate the residual field $\Delta \mathcal{V}_{t \rightarrow 1}^l$, which is added to the up-sampled $\mathcal{V}_{t \rightarrow 1}^{l-1}$ to yield the motion field $\mathcal{V}_{t \rightarrow 1}^l$.

As mentioned previously, the symmetric field $\mathcal{V}_{t \rightarrow 1}^S$ is estimated using the motion estimator of BMBC [34], which is at the quarter resolution. It is used as the up-sampled $\mathcal{V}_{t \rightarrow 1}^0$ at level $l = 1$ in Figure 5. Then, the refinement is performed for two levels, and the half resolution $\mathcal{V}_{t \rightarrow 1}^2$ becomes the final result $\mathcal{V}_{t \rightarrow 1}^A$. Since $\mathcal{V}_{t \rightarrow 0}^A$ and $\mathcal{V}_{t \rightarrow 1}^A$ are refined separately, they become asymmetric. As shown in Figure 1(a) and Figure 2(d), the asymmetric fields may represent motion information more faithfully than the symmetric ones.

Recently, Zhao *et al.* [43] improved the matching performance of PWC-Net with a learnable occlusion mask.

Similarly, we employ a mask in ABMR-Net to improve the refinement performance. The source frame \bar{I}_t is an approximation of I_t , so it may contain errors around motion boundaries or on complicated texture. Such errors make the matching process unreliable. Hence, we adopt the reliability mask Z_t^{l-1} in Figure 5 to suppress errors in anchor features. Notice that a source frame is masked in this work, while a target frame is masked in [43]. At level 1, O_t^0 is initialized to zero, and the up-sampled Z_t^0 is set to

$$\uparrow Z_t^0 = \exp\left(-\beta \cdot |\phi_B(\mathcal{V}_{t \rightarrow 0}^S, I_0) - \phi_B(\mathcal{V}_{t \rightarrow 1}^S, I_1)|\right) \quad (16)$$

where $\beta = 20$. If a certain pixel in \bar{I}_t has a large symmetric matching error in $|\phi_B(\mathcal{V}_{t \rightarrow 0}^S, I_0) - \phi_B(\mathcal{V}_{t \rightarrow 1}^S, I_1)|$, its feature is suppressed by the reliability mask.

3.2. Frame Synthesis

In Figure 4, we synthesize an intermediate frame I_t using two input frames I_0 and I_1 . We use four motion fields generated by the proposed ABME in Figure 3: two symmetric fields $\mathcal{V}_{t \rightarrow 0}^S$ and $\mathcal{V}_{t \rightarrow 1}^S$ and two asymmetric fields $\mathcal{V}_{t \rightarrow 0}^A$ and $\mathcal{V}_{t \rightarrow 1}^A$.

Using a feature extractor for frame synthesis, we extract multi-scale feature maps C_0^l from I_0 , where there are three levels $l \in \{1, 2, 3\}$. Note that the extraction from I_1 to C_1^l is performed similarly with shared parameters. The highest level maps C_0^3 and C_1^3 have the same spatial resolution as the input frames. We backward warp the input images and their feature pyramids. In Figure 4, I_0 is warped by $\mathcal{V}_{t \rightarrow 0}^S$ and $\mathcal{V}_{t \rightarrow 0}^A$ to obtain the estimate \hat{I}_t of I_t , while I_1 by $\mathcal{V}_{t \rightarrow 1}^S$ and $\mathcal{V}_{t \rightarrow 1}^A$. Thus, there are four candidate warped frames \hat{I}_t in total. Similarly, at each level l , there are four candidate warped features \hat{C}_t^l . These candidates are combined in a complementary manner to reconstruct the intermediate frame I_t more faithfully.

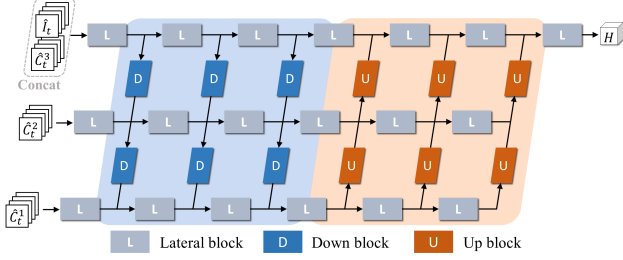


Figure 6. Modified GridNet for FilterNet

Given multiple candidates, we can synthesize I_t by blending them with weights. This simple blending, however, may cause blurry artifacts and reconstruction errors in occluded regions. Recent algorithms hence employ synthesis networks, which process warped frames to generate the intermediate frame in various ways: direct frame generation [30, 31], residual frame generation [2, 3], or dynamic local blending [34]. While the dynamic local blending synthesizes each pixel using local neighbors, the other two approaches use global contexts. To exploit both local and global information, we propose a novel synthesis network composed of two subnetworks: FilterNet and RefineNet.

FilterNet: It learns to generate dynamic filters for combining the four candidates, denoted by \hat{I}_t^c , $1 \leq c \leq 4$. We employ the modified version [30] of GridNet [12] as the backbone of FilterNet, which is shown in Figure 6. FilterNet takes the four candidates with the corresponding feature pyramids as input, by employing the lateral blocks. For each pixel (x, y) , the rightmost lateral block generates filter coefficients dynamically to fuse the 3×3 local neighboring pixels in each candidate. The coefficients are denoted by

$$H_{x,y}(i, j, c), \quad -1 \leq i, j \leq 1, \quad 1 \leq c \leq 4 \quad (17)$$

where (i, j) are local coordinates around (x, y) , and c is the candidate index. The filter coefficients are normalized, $\sum_c \sum_i \sum_j H_{x,y}(i, j, c) = 1$. Then, we obtain the filtered frame via the dynamic local convolution (DLC),

$$\tilde{I}_t(x, y) = \sum_{c=1}^4 \sum_{i=-1}^1 \sum_{j=-1}^1 H_{x,y}(i, j, c) \hat{I}_t^c(x + i, y + j). \quad (18)$$

Furthermore, by applying the same dynamic filters to the four warped feature candidates \hat{C}_t^l at the highest level $l = 3$, we obtain the filtered feature map \tilde{C}_t .

RefineNet: The dynamic filters consider local neighbors only. Thus, if the local neighbors do not contain proper information on a certain pixel due to motion errors or severe occlusion, its filtered result becomes also erroneous. To overcome this limitation using global information, RefineNet generates a residual frame ΔI_t to refine the filtered frame \tilde{I}_t . It has the same network architecture as FilterNet,

but it takes the filtered frame \tilde{I}_t , the filtered feature map \tilde{C}_t , and the warped feature pyramids as input. In particular, the feature map \tilde{C}_t , which is obtained using the same dynamic filters for \tilde{I}_t , meaningfully increases the refinement performance, as will be discussed in Section 4. After generating the residual ΔI_t , the final reconstructed frame is given by

$$I_t = \tilde{I}_t + \Delta I_t. \quad (19)$$

4. Experiments

4.1. Training

We first train the symmetric bilateral motion estimator and then, after fixing it, train ABMR-Net. Finally, we end-to-end train the frame synthesis network with those two networks.

ABME: We adopt the motion estimator of BMBC [34] for the symmetric bilateral motion estimation in Figure 3, but we retrain it by setting the stride of the first convolution layer to 2 for efficiency. To train the motion estimators of both BMBC and ABMR-Net, we define the photometric loss between a ground-truth I_t^{GT} and two warped frames as

$$\mathcal{L}_{\text{pho}} = \rho(I_t^{\text{GT}} - \phi_B(\mathcal{V}_{t \rightarrow 0}^A, I_0)) + \rho(I_t^{\text{GT}} - \phi_B(\mathcal{V}_{t \rightarrow 1}^A, I_1)) + \mathcal{L}_{\text{cen}}(I_t^{\text{GT}}, \phi_B(\mathcal{V}_{t \rightarrow 0}^A, I_0)) + \mathcal{L}_{\text{cen}}(I_t^{\text{GT}}, \phi_B(\mathcal{V}_{t \rightarrow 1}^A, I_1))$$

where $\rho(x) = (x^2 + \epsilon^2)^\alpha$ is the Charbonnier function [5] and \mathcal{L}_{cen} is the census loss [26, 44, 45] defined as the soft Hamming distance between census-transformed image patches [26] of size 7×7 . The parameters are set to $\alpha = 0.5$ and $\epsilon = 10^{-6}$.

To train the motion estimator of BMBC, we use the Adam optimizer [19] with a learning rate of $\eta = 10^{-4}$ until 0.1M iterations and halve η after every 0.04M iterations. We use a batch size of 24 for 0.2M iterations. For ABMR-Net, we also use the Adam optimizer with $\eta = 10^{-4}$ until 0.12M iterations and halve η after every 0.06M iterations. We use a batch size of 16 for 0.3M iterations.

Frame synthesis network: We define the synthesis loss \mathcal{L}_{syn} as the sum of the Charbonnier loss and the census loss between I_t^{GT} and its synthesized version I_t , given by

$$\mathcal{L}_{\text{syn}} = \rho(I_t^{\text{GT}} - I_t) + \mathcal{L}_{\text{cen}}(I_t^{\text{GT}}, I_t). \quad (20)$$

We use the Adam optimizer with $\eta = 10^{-4}$ until 0.35M iterations and halve η after every 0.15M iterations. We use a batch size of 6 for 0.8M iterations in total.

Training dataset: We use only the Vimeo90K training set [42] to train the proposed networks. It is composed of 51,312 triplets with a resolution of 448×256 . This training set does not contain motion ground-truth. We augment the dataset by randomly flipping, rotating, reversing the sequence order, and cropping 256×256 patches.

Table 1. Quantitative comparison (PSNR/SSIM) of video interpolation results on the UCF101, Vimeo90K, and SNU-FILM datasets. In each test, the best result is **boldfaced**, while the second best is underlined. All results are obtained by executing available source codes.

	UCF101	Vimeo90K	SNU-FILM				Runtime (seconds)	#Parameters (millions)
			Easy	Medium	Hard	Extreme		
ToFlow [42]	34.58/0.9667	33.73/0.9682	39.08/0.9890	34.39/0.9740	28.44/0.9180	23.39/0.8310	0.43	1.1
SepConv [33]	34.78/0.9669	33.79/0.9702	39.41/0.9900	34.97/0.9762	29.36/0.9253	24.31/0.8448	0.20	21.6
CyclicGen [22]	35.11/0.9684	32.09/0.9490	37.72/0.9840	32.47/0.9554	26.95/0.8871	22.70/0.8083	0.09	19.8
DAIN [2]	34.99/0.9683	34.71/0.9756	39.73/ 0.9902	35.46/ <u>0.9780</u>	<u>30.17/0.9335</u>	<u>25.09/0.8584</u>	0.13	24.0
CAIN [8]	34.91/ <u>0.9690</u>	34.65/0.9730	<u>39.89/0.9900</u>	<u>35.61/0.9776</u>	29.90/0.9292	24.78/0.8507	0.04	42.8
AdaCoF [20]	34.90/0.9680	34.47/0.9730	39.80/0.9900	35.05/0.9754	29.46/0.9244	24.31/0.8439	0.03	22.9
BMBC [34]	<u>35.15/0.9689</u>	<u>35.01/0.9764</u>	39.90/0.9902	35.31/0.9774	29.33/0.9270	23.92/0.8432	0.77	11.0
ABME (Proposed)	35.38/0.9698	36.18/0.9805	39.59/ <u>0.9901</u>	35.77/0.9789	30.58/0.9364	25.42/0.8639	0.22	18.1

4.2. Datasets

While we use strictly a single training dataset, we test the proposed ABME algorithm on various datasets.

UCF101 [38]: We use the test set constructed by Liu *et al.* [23], which contains 379 triplets of resolution 256×256 .

Vimeo90K [42]: The test set in Vimeo90K contains 3,782 triplets of spatial resolution 448×256 .

SNU-FILM [8]: It contains 1,240 triplets of videos of resolutions up to 1280×720 . It has four different settings – Easy, Medium, Hard, and Extreme.

Xiph [29]: It contains 30 raw video sequences for testing video codecs that have HD (1280×720) or FHD (1920×1080) resolutions. For triplets from the FHD sequences, we crop the central HD part of each frame without resizing. Thus, we extract 2,000 triplets of the HD resolution in total.

We divide these triplets into 5 classes (D1 ~ D5) according to the difficulty levels of interpolation. To quantify the difficulty, we compress each triplet using the recent video coding standard VVC [4] with a fixed QP. The middle frame is encoded in the B-frame mode, while the other two frames in the I-frame mode. The B-frame is motion-compensated from the others, and the motion vectors and compensation errors are encoded. The number of bits for the B-frame, hence, represents how difficult it is to interpolate it using the adjacent frames. The 2,000 triplets are sorted according to these numbers of bits and then grouped into the five classes so that each class contains 400 triplets. D1 is the easiest class, while D5 is the most difficult one.

X4K1000FPS [37]: Concurrently with this paper in ICCV 2021, Sim *et al.* present a high-quality, extensive dataset of 4K resolution. They provide the training set X-TRAIN and the test set X-TEST. We evaluate the proposed algorithm directly on X-TEST without retraining it on X-TRAIN.

4.3. Comparison with the State-of-the-Arts

We compare the proposed algorithm with conventional algorithms: ToFlow [42], SepConv [33], CyclicGen [22], DAIN [2], CAIN [8], AdaCoF [20], and BMBC [34].

Table 2. PSNRs on the Xiph dataset according to the difficulty levels. D1 is the easiest class, while D5 is the most difficult one.

	D1	D2	D3	D4	D5
DAIN [2]	34.65	<u>33.21</u>	<u>29.42</u>	<u>25.41</u>	22.61
CAIN [8]	<u>34.67</u>	32.68	27.97	24.98	<u>22.66</u>
AdaCoF [20]	34.23	32.16	27.53	24.84	22.28
BMBC [34]	34.47	32.13	27.21	24.52	22.39
ABME (Proposed)	35.21	34.15	30.26	25.77	23.02

Table 1 compares the average PSNR/SSIM scores. All results are obtained by executing available source codes.

- UCF101, Vimeo90K, and FILM (Medium) are easier to interpolate than FILM (Hard, Extreme). Also, on the easiest FILM (Easy), most algorithms interpolate high quality frames, and visual differences between the results are negligible.
- On all datasets except FILM (Easy), the proposed ABME algorithm provides the best results. Especially, on Vimeo90K, in comparison with the second best BMBC, ABME yields about 1 dB higher PSNR.
- On FILM (Medium, Hard, Extreme), DAIN achieves the second best results in general; on FILM (Hard), ABME yields 30.58 dB, while DAIN 30.17 dB.

In videos with faster motions, more regions are occluded around motion boundaries, making it more difficult to interpolate frames. Table 1 confirms that the proposed ABME handles these occluded regions effectively and provides excellent interpolation results.

Table 1 also lists the actual runtime for interpolating an intermediate frame in the ‘‘Urban’’ sequence in the Middlebury benchmark [1] using an RTX 2080 Ti GPU. The proposed algorithm is more than three times faster than BMBC.

Table 2 compares the average PSNRs on the Xiph dataset, which contains diverse sequences with challenging factors, such as complicated texture, fast motion, and severe occlusion. ABME outperforms all conventional algorithms in all difficulty classes. DAIN achieves the second-highest PSNRs in classes D2, D3, and D4. However, its motion and depth estimators were pre-trained using additional datasets, whereas ABME was trained strictly using Vimeo90K only. Nonetheless, in D2 and D3, ABME outperforms DAIN by



(a) GT (b) ToFlow (c) SepConv (d) CyclicGen (e) DAIN (f) CAIN (g) AdaCoF (h) BMBC (i) ABME

Figure 7. Qualitative comparison of interpolated frames. Triplets in SNU-FILM (Extreme) are used in this test. The proposed ABME algorithm in (i) reconstructs rapid objects faithfully to the ground truth in (a) without noticeable artifacts.

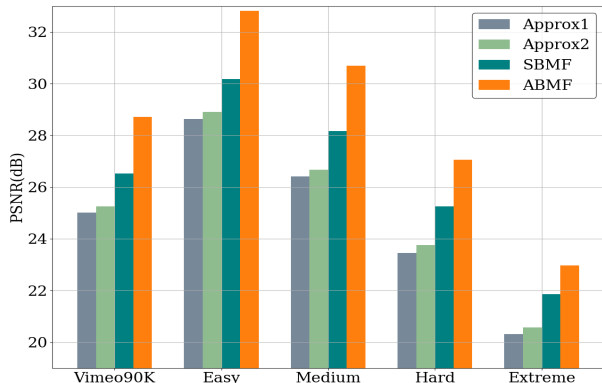


Figure 8. Comparison of different motion fields.

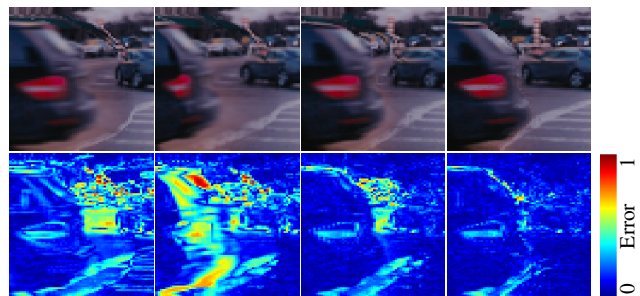
Table 3. Quantitative comparison on the X4K1000FPS dataset.

	PSNR	SSIM
DAIN [2]	26.78	0.8065
AdaCoF [20]	23.90	0.7271
XVFI [37]	<u>30.12</u>	<u>0.8704</u>
ABME (Proposed)	30.16	0.8793

0.94 dB and 0.84 dB, respectively. These results indicate that ABME interpolates challenging videos faithfully.

Table 3 compares the performances on the recent X4K1000FPS dataset. Its training set, X-TRAIN, is not used to retrain ABME. In contrast, X-TRAIN is used to train XVFI [37] that is designed for 4K sequences with extreme motions. Nevertheless, ABME performs slightly better than XVFI. It is a future research issue to tailor ABME for 4K sequences. For this purpose, the Vimeo90K training set might not be adequate, and X-TRAIN would be useful.

Figure 7 shows interpolation results of two frames in FILM (Extreme). Because of large movements and extreme deformation, the conventional algorithms fail to reconstruct details within the green and red squares reliably. In contrast, the proposed algorithm reconstructs them faithfully without any noticeable artifacts.



(a) Approx1 (b) Approx2 (c) SBMF (d) ABMF

Figure 9. Comparison of warped frames and their error maps.

4.4. Model Analysis

Let us analyze the contributions of key components in the proposed algorithm.

Motion fields: We compare the motion fields in Figure 2 quantitatively. Specifically, we test the approximate motion field in (7), called Approx1, and the approximate motion field in (9), called Approx2, and the symmetric bilateral motion field (SBMF), and the proposed asymmetric bilateral motion field (ABMF). PWC-Net [39] is used as the optical flow estimator for Approx1 and Approx2, while BMBC [34] for SBMF. These four types of motion fields are used, respectively, to warp input frames backward to approximate the intermediate frame, and then the PSNRs of the warped frames are computed. Figure 8 compares the average PSNRs on Vimeo90K and SNU-FILM. Approx1 produces the worst performance. Approx2, which improves the approximation via (9), slightly increases the PSNRs. SBMF in [34] yields considerably higher PSNRs than Approx1 or Approx2. However, the proposed ABMF outperforms SBMF significantly by more than 2 dB on Vimeo90K, Easy, and Medium and at least 1 dB on Hard and Extreme.

Figure 9 compares warped frames with error maps. The car moves quickly, causing severe occlusion. Hence, Ap-

Table 4. Impacts of the reliability mask Z_t and the offset map O_t on the asymmetric bilateral motion refinement. Average PSNRs of warped frames are reported.

Component		Vimeo90K	Extreme
Z_t	O_t	PSNR	PSNR
		27.81	22.75
✓		28.32	22.96
	✓	28.40	22.98
✓	✓	28.80	23.07

Table 5. Impacts of candidate warped frames on the frame synthesis. Average PSNRs of interpolated frames are reported.

Candidates	Vimeo90K	Extreme
	PSNR	PSNR
Symmetric only	35.91	25.28
Asymmetric only	36.05	25.34
Both	36.18	25.42

prox1 and Approx2 distort the warped frames severely. SBMF estimates the motion of the car accurately but causes errors around the boundary of the car due to disocclusion. In contrast, the proposed ABMF reduces these errors effectively by employing asymmetric motion vectors.

Reliability mask and offset map: We analyze the effectiveness of ABMR-Net in Figure 5. It uses a learnable reliability mask Z_t and a learnable offset map O_t . Table 4 shows that both components contribute to the overall performance of the asymmetric bilateral motion refinement.

Candidate warped frames: In the proposed frame synthesis in Figure 4, four candidate warped frames are used in the default mode. Specifically, two candidates are obtained using the symmetric fields $\mathcal{V}_{t \rightarrow 0}^S$ and $\mathcal{V}_{t \rightarrow 1}^S$, and the other two from the asymmetric fields $\mathcal{V}_{t \rightarrow 0}^A$ and $\mathcal{V}_{t \rightarrow 1}^A$. We test the following three combinations.

- Symmetric only: Two candidates using the symmetric fields are used.
- Asymmetric only: Two candidates using the asymmetric fields are used.
- Both: All four candidates are used.

Table 5 compares the average PSNRs of these three settings. First, ‘Symmetric only’ yields the worst performance, while it exceeds the performance of BMBC [34]. This means that the proposed synthesis network is more effective than that of BMBC. Second, ‘Asymmetric only’ performs better than ‘Symmetric only’ by loosening the linear motion constraint. Last, the best PSNRs are achieved by employing all four candidates, which complement each other to improve the interpolation performance.

Dynamic filters: Figure 10 shows the four candidate warped frames, their filter coefficient maps, and their error maps. To visualize dynamic filters, we compute the

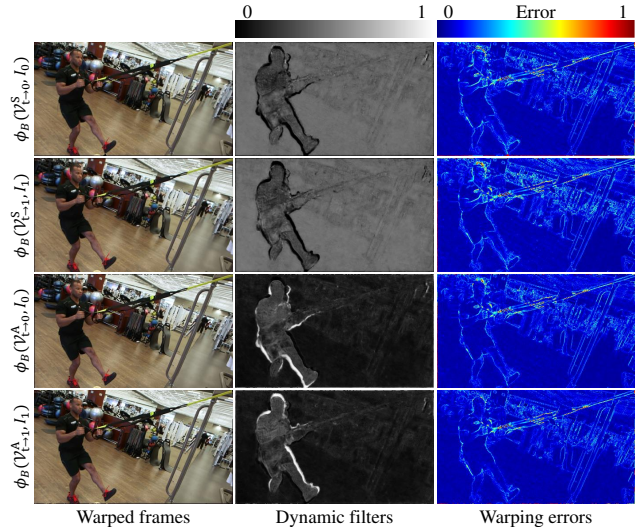


Figure 10. Visualization of dynamic filters.

absolute sum of coefficients for each pixel and render the sum in a gray level. The top two candidates are obtained by the symmetric fields, while the others by the asymmetric ones. The former are more reliable in the background without motion, while the latter are more effective around motion boundaries. Hence, the filter coefficients are determined accordingly. Moreover, to the left and right sides of the man, the third and fourth candidates are mainly used for the dynamic filtering, respectively. This is because different sides are occluded in different frames I_0 and I_1 . In such occluded regions, the linear motion constraint is invalid and the symmetric fields are not reliable.

Filtered feature maps: In Figure 4, RefineNet takes the filtered feature map \tilde{C}_t , as well as the filtered frame \tilde{I}_t , as input. If \tilde{C}_t is removed and only \tilde{I}_t is used to synthesize I_t , the average PSNR is reduced by 0.1 dB on Vimeo90K. Since \tilde{C}_t conveys contextual information in \tilde{I}_t , it helps RefineNet to remove noise in \tilde{I}_t and restore I_t more faithfully.

5. Conclusions

We proposed a novel, effective video frame interpolation algorithm. First, we developed the ABME technique to refine symmetric bilateral motions by loosening the linear motion constraint. Second, we designed the novel synthesis network that generates a set of dynamic filters and a residual frame using local and global information. Extensive experiments demonstrated that the proposed algorithm achieves state-of-the-art performance on various datasets.

Acknowledgements

This work was supported by the National Research Foundation of Korea (NRF) grants funded by the Korea government (MSIT) (Nos. NRF-2018R1A2B3003896, NRF-2019R1A2C4069806, and NRF-2021R1A4A1031864).

References

- [1] Simon Baker, Daniel Scharstein, J. P. Lewis, Stefan Roth, Michael J. Black, and Richard Szeliski. A database and evaluation methodology for optical flow. *Int. J. Comput. Vis.*, 92(1):1–31, Mar. 2011. **6**
- [2] Wenbo Bao, Wei-Sheng Lai, Chao Ma, Xiaoyun Zhang, Zhiyong Gao, and Ming-Hsuan Yang. Depth-aware video frame interpolation. In *CVPR*, pages 3703–3712, June 2019. **1, 2, 5, 6, 7**
- [3] 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 Trans. Pattern Anal. Mach. Intell.*, 43(3):933–948, Mar. 2021. **1, 2, 5**
- [4] Benjamin Bross, Jianle Chen, Jens-Rainer Ohm, Gary J. Sullivan, and Ye-Kui Wang. Developments in international video coding standardization after AVC, with an overview of Versatile Video Coding (VVC). *Proc. IEEE*, 109, 2021. **6**
- [5] Pierre Charbonnier, Laure Blanc-Feraud, Gilles Aubert, and Michel Barlaud. Two deterministic half-quadratic regularization algorithms for computed imaging. In *ICIP*, pages 168–172, Nov. 1994. **5**
- [6] Xianhang Cheng and Zhenzhong Chen. Video frame interpolation via deformable separable convolution. In *AAAI*, page 10607–10614, Feb. 2020. **1**
- [7] Byeong-Doo Choi, Jong-Woo Han, Chang-Su Kim, and Sung-Jea Ko. Motion-compensated frame interpolation using bilateral motion estimation and adaptive overlapped block motion compensation. *IEEE Trans. Circuit Syst. Video Technol.*, 17(4):407–416, Apr. 2007. **1**
- [8] 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, Feb. 2020. **1, 2, 6**
- [9] Alexey Dosovitskiy, Philipp Fischer, Eddy Ilg, Philip Hausser, Caner Hazirbas, Vladimir Golkov, Patrick Van Der Smagt, Daniel Cremers, and Thomas Brox. FlowNet: Learning optical flow with convolutional networks. In *ICCV*, pages 2758–2766, Dec. 2015. **1**
- [10] Karl M. Fant. A nonaliasing, real-time spatial transform technique. *IEEE Comput. Graph. Appl.*, 6(1):71–80, Jan. 1986. **2**
- [11] John Flynn, Ivan Neulander, James Philbin, and Noah Snavely. DeepStereo: Learning to predict new views from the world’s imagery. In *CVPR*, pages 5515–5524, June 2016. **1**
- [12] Damien Fourure, Rémi Emonet, Éliisa Fromont, Damien Muselet, Alain Treméau, and Christian Wolf. Residual conv-deconv grid network for semantic segmentation. In *BMVC*, pages 181.1–181.13, Sept. 2017. **5**
- [13] Shurui Gui, Chaoyue Wang, Qihua Chen, and Dacheng Tao. FeatureFlow: Robust video interpolation via structure-to-texture generation. In *CVPR*, pages 14004–14013, June 2020. **1**
- [14] 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 *CVPR*, pages 2462–2470, July 2017. **1**
- [15] Seong-Gyun Jeong, Chul Lee, and Chang-Su Kim. Motion-compensated frame interpolation based on multihypothesis motion estimation and texture optimization. *IEEE Trans. Image Process.*, 22(11):4497–4509, Nov. 2013. **1**
- [16] 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 *CVPR*, pages 9000–9008, June 2018. **1, 2, 3**
- [17] Rico Jonschkowski, Austin Stone, Jonathan T. Barron, Ariel Gordon, Kurt Konolige, and Anelia Angelova. What matters in unsupervised optical flow. In *ECCV*, pages 557–572, Aug. 2020. **1**
- [18] Nima Khademi Kalantari, Ting-Chun Wang, and Ravi Ramamoorthi. Learning-based view synthesis for light field cameras. *ACM Trans. Graph.*, 35(6):1–10, 2016. **1**
- [19] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR*, May 2015. **5**
- [20] Hyeongmin Lee, Taeoh Kim, Tae-young Chung, Daehyun Pak, Yuseok Ban, and Sangyoun Lee. AdaCoF: Adaptive collaboration of flows for video frame interpolation. In *CVPR*, pages 5316–5325, June 2020. **1, 2, 6, 7**
- [21] Liang Liu, Jiangning Zhang, Ruifei He, Yong Liu, Yabiao Wang, Ying Tai, Donghao Luo, Chengjie Wang, Jilin Li, and Feiyue Huang. Learning by analogy: Reliable supervision from transformations for unsupervised optical flow estimation. In *CVPR*, pages 6489–6498, June 2020. **1**
- [22] Yu-Lun Liu, Yi-Tung Liao, Yen-Yu Lin, and Yung-Yu Chuang. Deep video frame interpolation using cyclic frame generation. In *AAAI*, pages 8794–8802, Jan. 2019. **1, 2, 6**
- [23] Ziwei Liu, Raymond A Yeh, Xiaoou Tang, Yiming Liu, and Aseem Agarwala. Video frame synthesis using deep voxel flow. In *ICCV*, pages 4463–4471, Oct. 2017. **1, 2, 6**
- [24] Gucan Long, Laurent Kneip, Jose M Alvarez, Hongdong Li, Xiaohu Zhang, and Qifeng Yu. Learning image matching by simply watching video. In *ECCV*, pages 434–450, Oct. 2016. **1**
- [25] Guo Lu, Xiaoyun Zhang, Li Chen, and Zhiyong Gao. Novel integration of frame rate up conversion and HEVC coding based on rate-distortion optimization. *IEEE Trans. Image Process.*, 27(2):678–691, Feb. 2018. **1**
- [26] Simon Meister, Junhwa Hur, and Stefan Roth. UnFlow: Unsupervised learning of optical flow with a bidirectional census loss. In *AAAI*, Feb. 2018. **5**
- [27] Simone Meyer, Abdelaziz Djelouah, Brian McWilliams, Alexander Sorkine-Hornung, Markus Gross, and Christopher Schroers. PhaseNet for video frame interpolation. In *CVPR*, pages 498–507, June 2018. **1**
- [28] Simone Meyer, Oliver Wang, Henning Zimmer, Max Grosse, and Alexander Sorkine-Hornung. Phase-based frame interpolation for video. In *CVPR*, pages 1410–1418, June 2015. **1**
- [29] Christopher Montgomery. Xiph.org Video Test Media (derf’s collection), the Xiph Open Source Community. *Online*, <https://media.xiph.org/video/derf/>, 1994. **6**

- [30] Simon Niklaus and Feng Liu. Context-aware synthesis for video frame interpolation. In *CVPR*, pages 1701–1710, June 2018. [1](#), [2](#), [5](#)
- [31] Simon Niklaus and Feng Liu. Softmax splatting for video frame interpolation. In *CVPR*, pages 5437–5446, June 2020. [1](#), [2](#), [5](#)
- [32] Simon Niklaus, Long Mai, and Feng Liu. Video frame interpolation via adaptive convolution. In *CVPR*, pages 670–679, July 2017. [1](#)
- [33] Simon Niklaus, Long Mai, and Feng Liu. Video frame interpolation via adaptive separable convolution. In *ICCV*, pages 261–270, Oct. 2017. [1](#), [2](#), [6](#)
- [34] Junheum Park, Keunsoo Ko, Chul Lee, and Chang-Su Kim. BMBC: Bilateral motion estimation with bilateral cost volume for video interpolation. In *ECCV*, pages 109–125, Aug. 2020. [1](#), [2](#), [3](#), [4](#), [5](#), [6](#), [7](#), [8](#)
- [35] Anurag Ranjan and Michael J. Black. Optical flow estimation using a spatial pyramid network. In *CVPR*, pages 4161–4170, July 2017. [1](#)
- [36] 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 *ICCV*, pages 892–900, Oct. 2019. [2](#), [3](#)
- [37] Hyeonjun Sim, Jihyong Oh, and Munchurl Kim. XVFI: Extreme video frame interpolation. In *ICCV*, Oct. 2021. (arXiv preprint arXiv:2103.16206v2). [6](#), [7](#)
- [38] Khurram Soomro, Amir Roshan Zamir, and Mubarak Shah. UCF101: A dataset of 101 human actions classes from videos in the wild. *arXiv preprint arXiv:1212.0402*, 2012. [6](#)
- [39] Deqing Sun, Xiaodong Yang, Ming-Yu Liu, and Jan Kautz. PWC-Net: CNNs for optical flow using pyramid, warping, and cost volume. In *CVPR*, pages 8934–8943, June 2018. [1](#), [4](#), [7](#)
- [40] George Wolberg. *Digital Image Warping*. IEEE Computer Society Press, 1990. [2](#)
- [41] George Wolberg, H. M. Sueyllam, M. A. Ismail, and K. M. Ahmed. One-dimensional resampling with inverse and forward mapping functions. *J. Graph. Tools*, 5(3):11–33, 2000. [1](#)
- [42] Tianfan Xue, Baian Chen, Jiajun Wu, Donglai Wei, and William T. Freeman. Video enhancement with task-oriented flow. *Int. J. Comput. Vis.*, 127(8):1106–1125, Feb. 2019. [1](#), [2](#), [5](#), [6](#)
- [43] Shengyu Zhao, Yilun Sheng, Yue Dong, Eric I-Chao Chang, and Yan Xu. MaskFlowNet: Asymmetric feature matching with learnable occlusion mask. In *CVPR*, pages 6278–6287, June 2020. [1](#), [4](#)
- [44] Yiran Zhong, Pan Ji, Jianyuan Wang, Yuchao Dai, and Hongdong Li. Unsupervised deep epipolar flow for stationary or dynamic scenes. In *CVPR*, pages 12095–12104, June 2019. [5](#)
- [45] Yuliang Zou, Zelun Luo, and Jia-Bin Huang. DF-Net: Unsupervised joint learning of depth and flow using cross-task consistency. In *ECCV*, pages 36–53, Sept. 2018. [5](#)