# Dual Motion Attention and Enhanced Knowledge Distillation for Video Frame Interpolation

Dengyong Zhang*, Runqi Lou*, Xiaoping Hna§, Jiaxin Chen*, Xin Liao†, Gaobo Yang†, and Xiangling Ding‡

* School of Computer and Communication Engineering, Changsha University of Science and Technology, Changsha 410114, China

E-mail: Dengyong Zhang, zhdy@csust.edu.cn; Runqi Lou, lourunqi@stu.csust.edu.cn; Jiaxin Chen, jxchen@csust.edu.cn

† College of Computer Science and Electronic Engineering, Hunan University, Changsha 410082, China

E-mail: Xin Liao, xinliao@hnu.edu.cn; Gaobo Yang, yanggaobo@hnu.edu.cn

‡ School of Computer and Communication Engineering, Hunan University of Science and Technology, Xiangtan, 411201, Hunan, China

E-mail: xianglingding@163.com

§ Army Academy of Armored Forces, Beijing, 100072, China

E-mail: 15811060576@163.com

*Abstract*—**Video frame interpolation presents a formidable challenge in the domain of video generation, primarily due to the intricate motion dynamics exhibited by objects within video frames. With the advancements in deep learning, numerous flow-based methods for video frame interpolation have emerged. These methods aim to predict intermediate frames by leveraging the estimation of motion information between frames. In this paper, we propose a novel framework for modeling input video frames, which employs a coarse-to-fine structure to extract motion information between frames. Additionally, it incorporates a Bidirectional Correlation Volume and a complementary module of contextual features, specifically designed to pay attention to the symmetry of the optical flow and shallow motion features. We incorporate this dual attention to the knowledge distillation part of the model, which further improves the performance of the model. Leveraging this framework, our model demonstrates the ability to accurately predict motion information between frames, consequently producing visually appealing intermediate frames.The code is available at https://github.com/famt0531/DAEK.**

## I. Introduction

Video frame interpolation(VFI) serves as a low-level task in computer vision, utilized across various applications including video post-processing and surveillance. Its objective is to enhance the frame rate of a video sequence by generating intermediate frames between input frames. This process aims to achieve smoother video playback and mitigate motion blur[1], [2]. Traditional approaches typically involve estimating the optical flow between frames and subsequently interpolating or extrapolating along the optical flow vector. While effective with accurate optical flow, this method often yields significant artifacts and blurring when the optical flow estimation is imprecise.Moreover, they exhibit poor performance in scenarios involving occlusion and variations in luminance.

Recent years, numerous flow-based VFI techniques have emerged as prominent contenders in this domain[3]–[7]. These methods employ bidirectional optical flow to discern motion information across frames and map corresponding pixel positions. Utilizing estimated optical flow, they guide the warping process by transforming input images to interpolated frame positions, and blending them to preserve spatio-temporal correlations within the expected motion[4], [5]. Despite their promising outcomes, two main issues persist:

- Overlooking the correlation when estimating bidirectional optical flow, resulting in accumulated optical flow errors.
- Paying few attention to deep network features, which leads to missing detailed textures in the inference results.

To address the above limitations, in this paper, we propose a novel network architecture termed DAEK, which extracts coarse motion information, comprising intermediate optical flow and features, via three CNN-based Motion Estimation Blocks(MEB) in a coarse-to-fine fashion. Subsequently, these coarse motion information undergoes refinement through the Motion Refinement Blocks(MRB). During training, akin to RIFE, we employ intermediate supervision. However, we prioritize the accuracy of the generated optical flow of teacher model so that we use a combination of MEB and MRB, instead of directly employing the teacher module in RIFE. Our contributions are summarized as follows:

- We employed bidirectional correlation volumes to focus on the symmetry of motion at different stages of the network. Additionally, we utilized feature point multiplication to address the loss of motion information in the deeper layers of the network.
- We applied our proposed dual motion attention mechanism to a smaller model, which we trained as a teacher model for knowledge distillation. The teacher model with dual motion attention demonstrates a superior ability to guide the student model in accurately inferring optical flow information, compared to conventional teacher models
- We proposed a novel video frame interpolation framework, which adopts a coarse-to-fine structure to capture the motion information and considers the motion symmetry and texture loss through double motion focus.

## II. RELATED WORK

### A. VFI

VFI, a complex task in low-level computer vision, has seen the emergence of numerous innovative methods in recent years, broadly categorized into flow-based and non-optical flow-based approaches for motion simulation. In recent research endeavors, optical flow estimation has emerged as an indispensable component. For instance, Kong et al.[6] proposed IFRNet, a video interpolation network with a single encoder-decoder structure. It extracts feature pyramids from two input frames, refines bidirectional intermediate optical flow fields, and restores the desired output at the input resolution. Huang et al.[8] introduced RIFE, which estimates intermediate flows efficiently from coarse to fine, enhancing speed. They also developed a privileged distillation scheme for training for boosting performance. Jin et al.[9] created a compact model that estimates bidirectional motion simultaneously using a flexible pyramid recurrent framework, fine-tuning components within optical flow research for improved performance.

### B. Correlation Volume

Correlation volumes play a fundamental role in representing matching costs across various computer vision tasks. In the domain of VFI, several studies [9]–[11] have proposed construction schemes inspired by PWC-Net [12]. However, these schemes often focus solely on the local region's cost volume during matching, leading to inaccuracies when the region undergoes distortion. In contrast, AMT[13], building upon RAFT[14], transforms the unidirectional correlation volume into a bidirectional one, effectively capturing multiscale correspondence between frames. This compact global motion representation facilitates precise optical flow prediction, particularly for large motions. In this work, we adopt the design of AMT for the correlation volume construction.

### C. Self Attention

Given the remarkable success achieved by Transformer in natural language processing, there has been a surge in introducing Transformer to computer vision tasks, yielding promising results. In VFI, Liu et al. proposed ConvTransformer[15] to model long-term dependencies between 2D video frames in both temporal and spatial domains. Similarly, Lu et al.[16] utilized Transformer as an encoder in their model, enabling it to learn relationships between each frame and others via a self-attentive mechanism. However, many VFI methods based on Transformer employ the self-attention mechanism as a feature extractor, which can be computationally expensive due to the large number of parameters involved. Therefore, rather than directly using the Transformer for feature extraction, we simulate its self-attention mechanism to supplement contextual information during optical flow updating.

## III. PROPOSED METHOD

### A. Problem Description

Given two frames $I^0$ and $I^1$, the objective of Video Frame Interpolation (VFI) is to generate the intermediate frame $I^t$,

where $t \in [0, 1]$ represents the time step and is typically set to 0.5 by default. Our proposed model follows an end-to-end optical-flow based approach, illustrated in Fig. 1. We use the same residual structure to extract appearance features and context features, which are used to construct bidirectional cost volumes and enrich the texture information in the deeper layers of the network, respectively. The coarse motion information is then refined in MRB.

### B. MEB

MEB iterates to estimate intermediate flows utilizing a coarse-to-fine structure, enabling effective handling of large motions. For each MEB, the inputs consist of $I^0$ and $I^1$ which obtained after bilinear interpolation downsampling, intermediate flows $flow_i \in \mathbb{R}^{4 \times H' \times W'}$, mask $M_i \in \mathbb{R}^{1 \times H' \times W'}$, intermediate features $Feat_i \in \mathbb{R}^{C' \times H' \times W'}$ (where $i = 1, 2$) which is predicted by the previous MEB, and the $\tilde{I}_0, \tilde{I}_1$ obtained by warping $I^0$ and $I^1$ with $flow_i$ (for the first MEB, we merely input $I^0$ and $I^1$). Specifically, we employ three consecutive motion estimation modules, each processing motion information at different scales. Following the multi-scale strategy outlined in RIFE[8], we set the scales to $[1/8, 1/4, 1/2]$.

All MEBs share the same structure, realized entirely by convolutional layers. Additionally, we output intermediate features for each MEB, as shown in Fig. 2, which play a crucial role in the subsequent optical flow update phase and motion estimation phase. By including intermediate features, our model retains more motion information during the iterative process of generating optical flow at multiple scales.
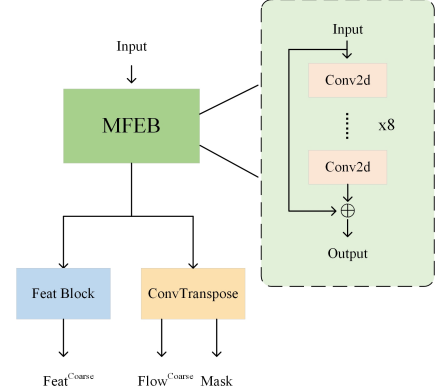


Fig. 2. The structure of MEB. The $Conv$ comprises 8 layers of convolutional layers and PRELU units. Each convolutional layer utilizes a 3x3 convolutional kernel with a stride of 1 and padding of 1.

### C. MRB

After obtaining the initial bi-directional optical flow and intermediate features, we utilize MRB to generate residuals for the bi-directional optical flow and intermediate features. These residuals are then used to update both the optical flow and intermediate features. We use bidirectional cost volumes to obtain the bidirectional matching scores between the two real frames. Our feature matching is consistently performed at
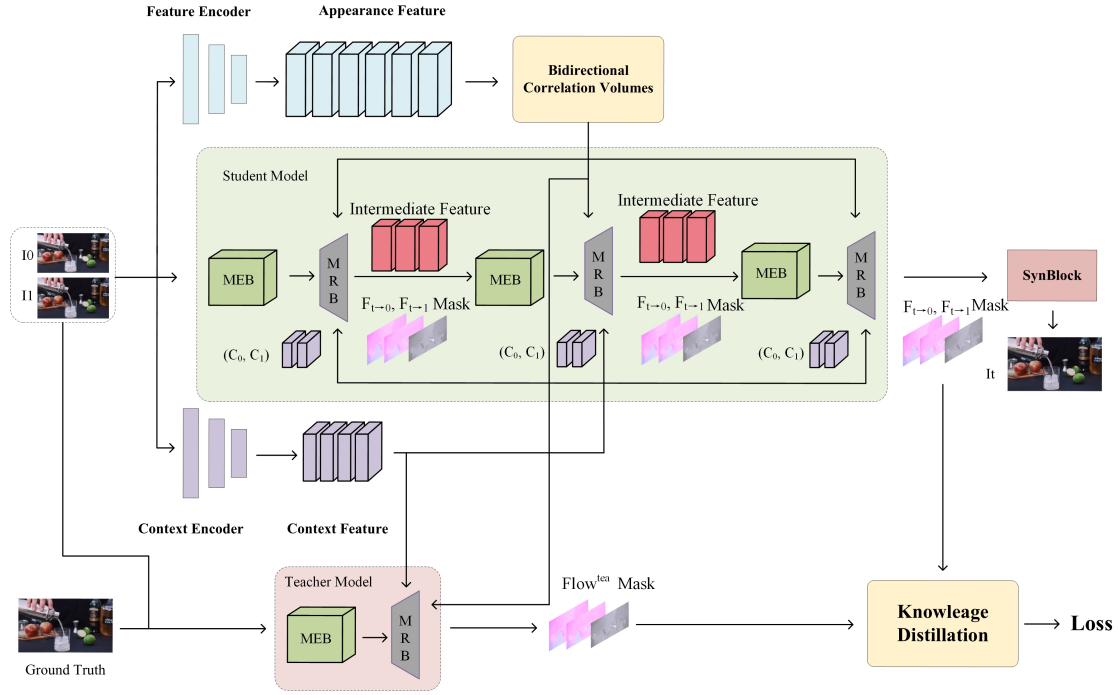
Fig. 1. The overview of the proposed model. The input frames are processed through MEB and MRB across three scales to generate bidirectional optical flow ($F_{t\to0}$, $F_{t\to1}$) and $Mask$. The synthesis module then outputs the intermediate frame. Our teacher module generates distillation loss only during the training phase and is discarded during inference.

a resolution of $\frac{1}{8}$. For each scale output of the MEB, we first downsample them to $\frac{1}{8}$ resolution using linear interpolation. After updates, we then restore them to the original resolution. We employ several convolutional layers to extract bidirectional correlation features $BCF \in \mathbb{R}^{C \times \frac{H}{8} \times \frac{W}{8}}$ and coarse flow features, then concatenate them along the channel dimension to obtain the coarse motion information $CMF$, as follows:

$$CMF = Cat(Conv(Flow^{coarse}), Conv(BCF)) \quad (1)$$

Where $Cat$ represents the concat operation, and $Conv$ for convolution operation.

We posit that features at the shallow end of the network encompass more textural details, whereas deeper layers tend to encapsulate semantic information. However, as the network depth increases, texture features become inevitably attenuated. Therefore, we introduce Contextual Feature Complementary Module(CFCM) that uses an attention mechanism to complement feature information at deeper layers of the network which serves as a crucial component of MRB.

As shown in Fig.3, the context feature undergoes transformation by a linear layer into $Query, Key \in \mathbb{R}^{C \times H*W}$. Subsequently, the resulting outputs are multiplied and passed through a $Softmax$ unit to obtain $Score$. Meanwhile, the $CMF$ is transformed into $Value \in \mathbb{R}^{C^v \times H^v \times W^v}$ through a convolutional layer with a kernel size of 1×1. The CFCM we propose is essentially a feature enhancement method, so we abandon the tedious calculation of multiplying the result of Q and K and then multiplying it with V, and instead add it to V. In order to better achieved feature fusion and feature
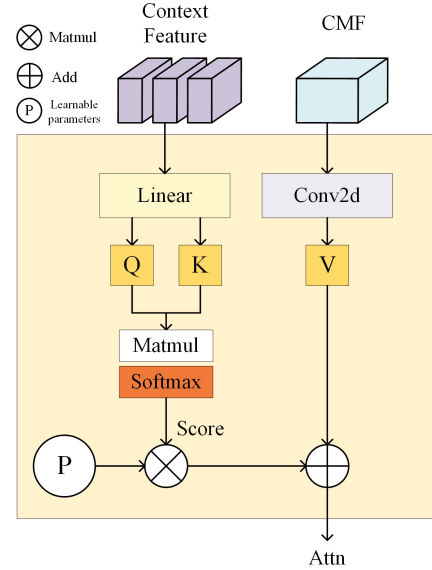


Fig. 3. The structure of CFCM, where the convolutional layer has a kernel size of 1x1.

enhancement, we introduce a learnable parameter, initialized to 0, which is multiplied with the Score and added to the Value to obtain Attn, This step can be expressed as:

$$Q, K = linear(context feature) \quad (2)$$

$$Score = Softmax(Q \odot K) \quad (3)$$

3

$$V = Conv(CMF) \tag{4}$$

$$Attn = V + (Para \odot Score) \tag{5}$$

Where $linear$ represents fully connected layer, $Softmax$ represents Softmax unit, $Conv$ represents convolution operation, $\odot$ represents multiplication.

Finally, we concatenate $Attn$ with $Flow^{coarse}$ and intermediate feature along the channel dimension and apply convolutional layers to obtain the residuals which is added with coarse flow and intermediate feature to achieve the goal of refining.

### D. Enhanced Knowleage Distillation

We stack an extra MEB as the optical flow estimation module for the teacher model, the MEB for the teacher model has the same structure as the MEB for the student model, except that extra ground truth is added as an input. In order to obtain a more accurate labeling of the optical flow, we add a corresponding MRB for the teacher model, which produces residuals only for the flow involved in knowledge distillation, but not for the intermediate feature, since the intermediate feature here are not involved in subsequent parameter updates during the training phase. After obtaining the optical flow labels generated by the teacher model again, we use it and the intermediate flows generated by the student model at different scales to generate distillation losses, as follows:

$$L_{dis} = \sum ||Flow_{t \to 0}^{tea}, Flow_{t \to 0}||_1 + ||Flow_{t \to 1}^{tea}, Flow_{t \to 1}||_1 \tag{6}$$

Where $Flow_{t \to i}^{tea}(i = 0, 1)$ represents flow generated by teacher module, $Flow_{t \to i}(i = 0, 1)$ represents flow fields generated by student module, and $|| \cdot ||_1$ is the operation of $L1$ loss.

## IV. EXPERIMENTS

### A. Datasets

*1) Training Dataset:* We train our DAEK on the Vimeo-90K[17] training set which comprises a total of 51,312 triples. Each triple contains three consecutive video frames with a resolution of 448 × 256. In training stage, we crop patches from the original images with a resolution of 224 × 224 and randomly augment it by performing horizontal and vertical flipping, chronological flipping, and rotating by 90 degrees, similar to other methods such as RIFE[8], AdaCoF[18].

*2) Evaluation settings:* We test the performance of DAEK on the Vimeo-90K test set, UCF101[19] dataset, and SNU-FILM[20] dataset, and for all the tables in this section, the best performance is labeled by us in bold font.

### B. Implementation Details

*1) Loss function:* DAEK uses three loss functions, which including reconstruction loss $L_{rec}$, teacher reconstruction loss $L_{rec}^{tea}$ and distillation loss $L_{dis}$. The total loss $L$:

$$L = L_{rec} + L_{rec}^{tea} + \lambda L_{dis} \tag{7}$$

where we set $\lambda$=0.01.

Our reconstruction loss quantifies the quality of the reconstructed $I_t$ and $I_t^{gt}$ using the $L1$ loss function with the Laplacian operator[21] which could be defined as:

$$L_{rec} = Lap(I_t, I_{gt}) \tag{8}$$

$$L_{rec}^{tea} = Lap(I_t^{tea}, I_{gt}) \tag{9}$$

where $Lap$ is the Laplacian operator, $I_t$ is the intermediate frame of student module, $I_{gt}$ is the ground truth, and $I_t^{tea}$ is the intermediate frame of teacher module.

*2) Training Details:* DAEK was implemented using the Pytorch framework and trained on a NVIDIA A30 Tensor Core GPU. We used AdamW[22] as the optimizer with a weight decay of $10^{-3}$. The initial learning rate is $10^{-4}$, which is decayed to $10^{-6}$ by cosine annealing strategy during the training process, and the whole training process requires 300 epochs with a batch size set to 32.

### C. Performance Evaluation

*1) Objective results:* We compare DAEK to current state-of-the-art methods including, IFRNet[23], AdaCoF[18], BMBC[11], EBME[9], RIFE[8], M2M[24], as shown in TabI.

We conducted a subjective performance comparison on the highly challenging extreme dataset of SNU-FILM, as shown in Fig.4 and Fig.5. When interpolating across multiple video frames, the motion within the frames becomes significantly more extensive and highly nonlinear. DAEK demonstrates excellent performance when dealing with these complex motions. Other optical flow-based methods exhibit noticeable blurring and distortion, while the kernel-based method AdaCoF, although maintaining motion continuity, suffers from pixel loss. In contrast, DAEK preserves rich texture details and accurately estimates motion continuity.

### D. Ablation Study

*1) Structural Ablation Experiment:* We use models without dual motion attention as the baseline, as shown in TabII. Here, $F$ indicates the use of intermediate features, while its absence indicates no use. $CA$ represents correlation attention, $SFA$ stands for shallow feature attention, and $DA$ denotes dual attention. The use of these two motion attention individually brings certain performance improvements. When used together, the model's ability to represent bidirectional motion is further enhanced. Additionally, the introduction of intermediate features provides the intermediate frames with realistic texture details and accurate motion representation.

TABLE II
VALIDITY OF DIFFERENT COMPONENTS OF OUR MODEL

| Setting | Vimeo90K PSNR | UCF01 PSNR |
|---|---|---|
| Baseline | 35.32 | 35.29 |
| DAEK-F-CA | 35.43 | 35.39 |
| DAEK-F-SFA | 35.31 | 35.35 |
| DAEK-F-DA | **35.66** | **35.45** |
| DAEK-DA | 35.49 | 35.30 |

TABLE I
THE EVALUATION OF VARIOUS INTERPOLATION METHODS ON VIMEO-90K TEST SET, UCF101 DATASET AND SNU-FILM DATASE.

| Method | Vimeo-90K | UCF101 | SNU-FILM | | | | Parameters (M) | Runtime (S) |
|--------|-----------|--------|------|--------|------|---------|----------------|-------------|
| | | | Easy | Medium | Hard | Extreme | | |
| IFRNet | **35.80**/**0.979** | 35.29/0.969 | **40.03**/0.990 | **35.94**/0.979 | 30.41/0.935 | 25.05/0.858 | **5.0** | 0.02 |
| AdaCoF | 32.00/0.971 | 35.08/0.966 | 39.80/0.990 | 35.05/0.975 | 29.46/0.924 | 24.31/0.843 | 21.8 | 0.02 |
| BMBC | 35.01/0.976 | 35.15/0.969 | 39.90/0.990 | 35.31/0.977 | 29.33/0.927 | 23.92/0.843 | 11.0 | 0.3 |
| EBME | 35.58/0.978 | 35.30/0.969 | 40.01/**0.991** | 35.80/0.979 | 30.42/0.935 | **25.25**/**0.861** | 3.9 | 0.02 |
| RIFE | 35.32/0.976 | 35.29/**0.973** | 38.06/0.983 | 35.14/0.975 | 29.72/0.925 | 24.31/0.843 | 9.8 | **0.01** |
| M2M-PWC | 35.49/0.978 | 35.32/0.970 | 39.66/**0.991** | 35.74/**0.980** | 30.32/**0.936** | 25.07/0.860 | 7.6 | 0.03 |
| Ours | 35.66/**0.979** | **35.42**/0.971 | 36.13/0.981 | 33.97/0.971 | 29.91/0.933 | 25.13/**0.861** | 22.48 | 0.03 |



(a)Ground Truth  (b)AdaCoF  (c)IFRNet  (d)M2M  (e)RIFE  (f)DAEK

Fig. 4.   Comparison of Visualized Test Results on Extreme set of SNU-FILM



(a)Ground Truth  (b)AdaCoF  (c)IFRNet  (d)M2M  (e)RIFE  (f)DAEK

Fig. 5.   Comparison of Visualized Test Results on Extreme set of SNU-FILM

*2) Knowledge Distillation Ablation Experiment:* Knowledge distillation can refine model parameters during the training phase, guiding it to produce more accurate results. As shown in Tab.III, the model's performance significantly declines without knowledge distillation. However, when we apply the knowledge distillation proposed in RIFE, the model's performance improves. Since we incorporate dual motion attention in the knowledge distillation process, the flow labels generated by the teacher model include richer bidirectional motion information, thereby guiding the student model to produce more accurate outputs. To validate the robustness of our enhanced knowledge distillation, we replaced the teacher model in RIFE. This substitution resulted in a performance improvement for RIFE as well.

TABLE III
ABLATION EXPERIMENTS FOR DISTILLATION SCHEMES.

| Setting | Vimeo90K |
|---------|----------|
| | PSNR |
| DAEK w/o distill | 35.20 |
| DAEK w/ RT | 35.40 |
| DAEK w/ OT | **35.66** |
| RIFE w/ OT | 35.49 |

## V. CONCLUSIONS

In this study, we introduce a novel model named DAEK, which initially generates the base motion flow field using MEB. Subsequently, the original motion flow field is refined using MRB which integrates dual motion attention. These designs aids DAEK in addressing complex motions in video

frames resulting in significant results. Moving forward, we aim to optimize the construction of bidirectional correlation volumes, for instance, by employing the PatchMatch[25], to reduce the computational load and inference time of the model. Additionally, we plan to enhance the model's architecture to improve its learning capacity, as a higher initial learning rate facilitates parameter optimization, ultimately boosting the model's learning effectiveness.

## VI. Acknowledgements

## References

[1] J. He, G. Yang, X. Liu, and X. Ding, "Spatio-temporal saliency-based motion vector refinement for frame rate up-conversion," *ACM Transactions on Multimedia Computing, Communications, and Applications (TOMM)*, vol. 16, no. 2, pp. 1–18, 2020.

[2] W. Bao, W. S. Lai, C. Ma, X. Zhang, Z. Gao, and M. H. Yang, "Depth-aware video frame interpolation," *IEEE*, 2019.

[3] Z. Zhang, L. Song, R. Xie, and L. Chen, "Video frame interpolation using recurrent convolutional layers," in *2018 IEEE Fourth International Conference on Multimedia Big Data (BigMM)*, IEEE, 2018, pp. 1–6.

[4] S. Yu, B. Park, and J. Jeong, "Posnet: 4x video frame interpolation using position-specific flow," *IEEE*, 2019.

[5] H. Li, Y. Yuan, and Q. Wang, "Fi-net: A lightweight video frame interpolation network using feature-level flow," *IEEE Access*, vol. PP, no. 99, pp. 1–1, 2019.

[6] L. Kong, B. Jiang, D. Luo, *et al.*, "Ifrnet: Intermediate feature refine network for efficient frame interpolation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 1969–1978.

[7] S. Niklaus and F. Liu, "Softmax splatting for video frame interpolation," *IEEE*, 2020.

[8] Z. Huang, T. Zhang, W. Heng, B. Shi, and S. Zhou, "Real-time intermediate flow estimation for video frame interpolation," in *European Conference on Computer Vision*, Springer, 2022, pp. 624–642.

[9] X. Jin, L. Wu, G. Shen, *et al.*, "Enhanced bi-directional motion estimation for video frame interpolation," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2023, pp. 5049–5057.

[10] Q. Hou, A. Ghildyal, and F. Liu, "A perceptual quality metric for video frame interpolation," in *European Conference on Computer Vision*, Springer, 2022, pp. 234–253.

[11] J. Park, K. Ko, C. Lee, and C.-S. Kim, "Bmbc: Bilateral motion estimation with bilateral cost volume for video interpolation," in *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XIV 16*, Springer, 2020, pp. 109–125.

[12] D. Sun, X. Yang, M. Y. Liu, and J. Kautz, "Pwc-net: Cnns for optical flow using pyramid, warping, and cost volume," in *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.

[13] Z. Li, Z.-L. Zhu, L.-H. Han, Q. Hou, C.-L. Guo, and M.-M. Cheng, "Amt: All-pairs multi-field transforms for efficient frame interpolation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 9801–9810.

[14] Z. Teed and J. Deng, "Raft: Recurrent all-pairs field transforms for optical flow," in *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part II 16*, Springer, 2020, pp. 402–419.

[15] Z. Liu, S. Luo, W. Li, *et al.*, "Convtransformer: A convolutional transformer network for video frame synthesis," *arXiv preprint arXiv:2011.10185*, 2020.

[16] L. Lu, R. Wu, H. Lin, J. Lu, and J. Jia, "Video frame interpolation with transformer," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 3532–3542.

[17] T. Xue, B. Chen, J. Wu, D. Wei, and W. T. Freeman, "Video enhancement with task-oriented flow," *International Journal of Computer Vision*, vol. 127, pp. 1106–1125, 2019.

[18] H. Lee, T. Kim, T. Y. Chung, D. Pak, and S. Lee, "Adacof: Adaptive collaboration of flows for video frame interpolation," in *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.

[19] K. Soomro, A. R. Zamir, and M. Shah, "Ucf101: A dataset of 101 human actions classes from videos in the wild," *Computer Science*, 2012.

[20] M. Choi, H. Kim, B. Han, N. Xu, and K. M. Lee, "Channel attention is all you need for video frame interpolation," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, 2020, pp. 10 663–10 671.

[21] S. Niklaus and F. Liu, "Context-aware synthesis for video frame interpolation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 1701–1710.

[22] D. Kingma and J. Ba, "Adam: A method for stochastic optimization," *Computer Science*, 2014.

[23] L. Kong, B. Jiang, D. Luo, *et al.*, "Ifrnet: Intermediate feature refine network for efficient frame interpolation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022.

[24] P. Hu, S. Niklaus, L. Zhang, S. Sclaroff, and K. Saenko, "Video frame interpolation with many-to-many splatting and spatial selective refinement," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023.

[25] Z. Zheng, N. Nie, Z. Ling, *et al.*, "Dip: Deep inverse patchmatch for high-resolution optical flow," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 8925–8934.