

AMT: All-Pairs Multi-Field Transforms for Efficient Frame Interpolation

Zhen Li* Zuo-Liang Zhu* Ling-Hao Han Qibin Hou Chun-Le Guo[†] Ming-Ming Cheng
VCIP, CS, Nankai University

{zhenli1031, nkuzhuzl}@gmail.com, lhhan@mail.nankai.edu.cn
{houqb, guochunle, cmm}@nankai.edu.cn

Abstract

We present *All-Pairs Multi-Field Transforms (AMT)*, a new network architecture for video frame interpolation. It is based on two essential designs. First, we build bidirectional correlation volumes for all pairs of pixels, and use the predicted bilateral flows to retrieve correlations for updating both flows and the interpolated content feature. Second, we derive multiple groups of fine-grained flow fields from one pair of updated coarse flows for performing backward warping on the input frames separately. Combining these two designs enables us to generate promising task-oriented flows and reduce the difficulties in modeling large motions and handling occluded areas during frame interpolation. These qualities promote our model to achieve state-of-the-art performance on various benchmarks with high efficiency. Moreover, our convolution-based model competes favorably compared to Transformer-based models in terms of accuracy and efficiency. Our code is available at <https://github.com/MCG-NKU/AMT>.

1. Introduction

Video frame interpolation (VFI) is a long-standing video processing technology, aiming to increase the temporal resolution of the input video by synthesizing intermediate frames from the reference ones. It has been applied to various downstream tasks, including slow-motion generation [23, 68], novel view synthesis [11, 31, 79], video compression [67], text-to-video generation [56], etc.

Recently, flow-based VFI methods [18, 23, 28, 37, 57, 77] have been predominant in referenced research due to their effectiveness. A common flow-based technique estimates bilateral/bidirectional flows from the given frames and then propagates pixels/features to the target time step via backward [2, 18, 28] or forward [15, 41, 42] warping. Thus, the quality of a synthesized frame relies heavily on flow estimation results. In fact, it is cumbersome to approximate intermediate flows through pretrained optical flow models,

*Equal contribution

[†]C.L. Guo is the corresponding author.

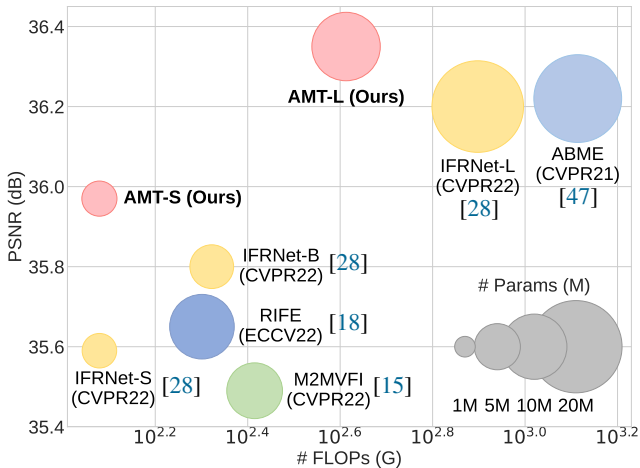


Figure 1. Performance vs. number of parameters and FLOPs. The PSNR values are obtained from the Vimeo90K dataset [73]. We use a 720p frame pair to calculate FLOPs. Our AMT outperforms the state-of-the-art methods and is with higher efficiency.

and these flows are unqualified for VFI usage [15, 18].

A feasible way to alleviate this issue is to estimate *task-oriented flows* in an end-to-end training manner [23, 28, 35, 73]. However, some major challenges, such as large motions and occlusions, are still pending to be resolved. These challenges mainly arise from the defective estimation of optical flows. Thus, a straightforward question should be: Why do previous methods have difficulties in predicting promising task-oriented flows when facing these challenges? Inspired by the recent studies [28, 73] that demonstrate *task-oriented flow is generally consistent with ground truth optical flow but diverse in local details*, we attempt to answer the above question from two perspectives:

(i) The flow fields predicted by existing VFI methods are *not consistent enough* with the true displacements, especially when encountering large motions (see Fig. 2). Existing methods mostly adopt the UNet-like architecture [52] with plain convolutions to build VFI models. However, this type of architecture is vulnerable to accumulating errors at early stages when modeling large motions [49, 61, 71, 78]. As a result, the predicted flow fields are not accurate.

(ii) Existing methods predict one pair of flow fields, re-

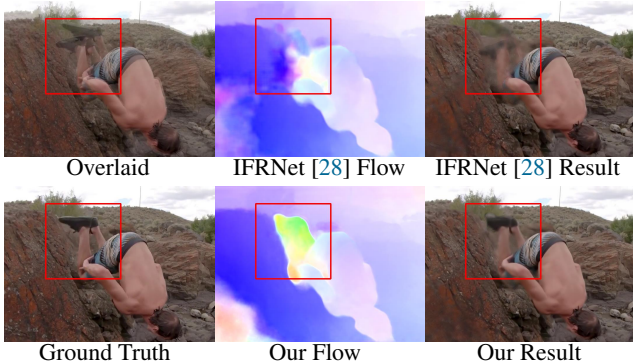


Figure 2. Qualitative comparisons of estimated flows and the interpolated frames. Our AMT guarantees the general consistency of intermediate flows and synthesizes fast-moving objects with occluded regions precisely, while the previous state-of-the-art IFRNet [28] fails to achieve them.

stricting the solution set in a tight space. This makes them struggle to handle occlusions and details around the motion boundaries, which consequently deteriorates the final results (see Fig. 2 and Fig. 5).

In this paper, we present a new network architecture, dubbed **All-pairs Multi-field Transforms (AMT)**, for video frame interpolation. AMT explores two new designs to improve the fidelity and diversity of predicted flows regarding the above two main shortcomings of previous works.

Our first design is based on all-pairs correlation in RAFT [61], which adequately models the dense correspondence between frames, especially for large motions. We propose to build *bidirectional correlation volumes* instead of a unidirectional one and introduce a *scaled* lookup strategy to solve the coordinate mismatch issue caused by the invisible frame. Besides, the retrieved correlations assist our model in *jointly* updating bilateral flows and the interpolated content feature in a *cross-scale* manner. Thus, the network guarantees the fidelity of flows across scales, laying the foundation for the following refinement.

Second, considering that predicting one pair of flow fields is hard to cope with the occlusions, we propose to derive multiple groups of fine-grained flow fields from one pair of updated coarse bilateral flows. The input frames can be separately backward warped to the target time step by these flows. Such diverse flow fields provide adequate potential solutions for each pixel to be interpolated, particularly alleviating the ambiguity issue in the occluded areas.

We examine the proposed AMT on several public benchmarks with different model scales, showing strong performance and high efficiency in contrast to the state-of-the-art (SOTA) methods (see Fig. 1). Our small model outperforms IFRNet-B, a SOTA lightweight model, by +0.17dB PSNR on Vimeo90K [73] with only 60% of its FLOPs and parameters. For the large-scale setting, our AMT exceeds the previous SOTA (*i.e.*, IFRNet-L) by +0.15 dB PSNR on Vimeo90K [73] with 75% of its FLOPs and 65% of its pa-

rameters. Besides, we provide a huge model for comparison with the SOTA transformer-based method VFIFormer [37]. Our convolution-based AMT shows a comparable performance but only needs nearly $23\times$ less computational cost compared to VFIFormer [37]. Considering its effectiveness, we hope our AMT could bring a new perspective for the architecture design in efficient frame interpolation.

2. Related Work

Video Frame Interpolation: The development of deep learning has spawned a large amount of VFI methods. These methods can be roughly divided into three categories: kernel-based [5, 6, 29, 43, 44, 48], hallucination-based [8, 25, 27, 54], and flow-based ones [1, 2, 18, 28, 35, 41, 42, 72, 73].

Kernel-based methods attempt to capture motion with dynamic kernel weights [43, 44, 48] or/and offsets [5, 6, 10, 29]. With the help of off-the-shelf architectures [9, 14, 53, 63], hallucination-based methods directly generate the interpolated frame from features of input pairs. Thanks to the robustness of optical flow, flow-based methods have become mainstream in VFI. Previous methods resort to a pre-trained flow model [41, 72] or a jointly trained estimation module [28, 35, 73] to obtain the flow estimation. For generating task-oriented flows, some methods [18, 28] propose intermediate supervisions to distill motion knowledge from the pseudo ground truth. Subsequently, backward warping [2, 18, 28] and forward warping [41, 42, 45] are standard schemes in the usage of estimated flows. A UNet-like architecture is a common choice [1, 2, 41, 42] to obtain the final synthesized frame, and the transformer [34, 65], as a prevailing architecture, is introduced [37, 54, 76] for a better synthesis. Recent works [28, 50] discard an independent synthetic network in consideration of efficiency. However, these methods suffer from the inability in modeling large motions and in dealing with occlusions.

Task-Oriented Flow: Initially, flow-based video processing methods [1, 23, 46, 72] estimate flows and process images individually. However, this two-step pipeline ignores the gap between true optical flow with task-specific objectives, which could be suboptimal for a specific task. ToFlow [73] proposes the concept of task-oriented flow, facilitating the development of video processing methods [3, 16, 30, 32, 74] significantly. Typically, the VFI-oriented flow is generally consistent with the true flow while diversifies in detail (*e.g.*, occluded regions). Super Slomo [23] introduces a mask to handle the occlusion explicitly and provides a standard formulation for synthesizing intermediate frames, which utilizes by following methods [7, 18, 28, 55] up to now. IFRNet [28] and RIFE [18] propose task-oriented flow distillation losses to provide a prior of intermediate flow in training. Different from them, we consider the estimation of task-oriented flows from *the perspective of architecture design*. We introduce all-pairs correlation to strengthen the

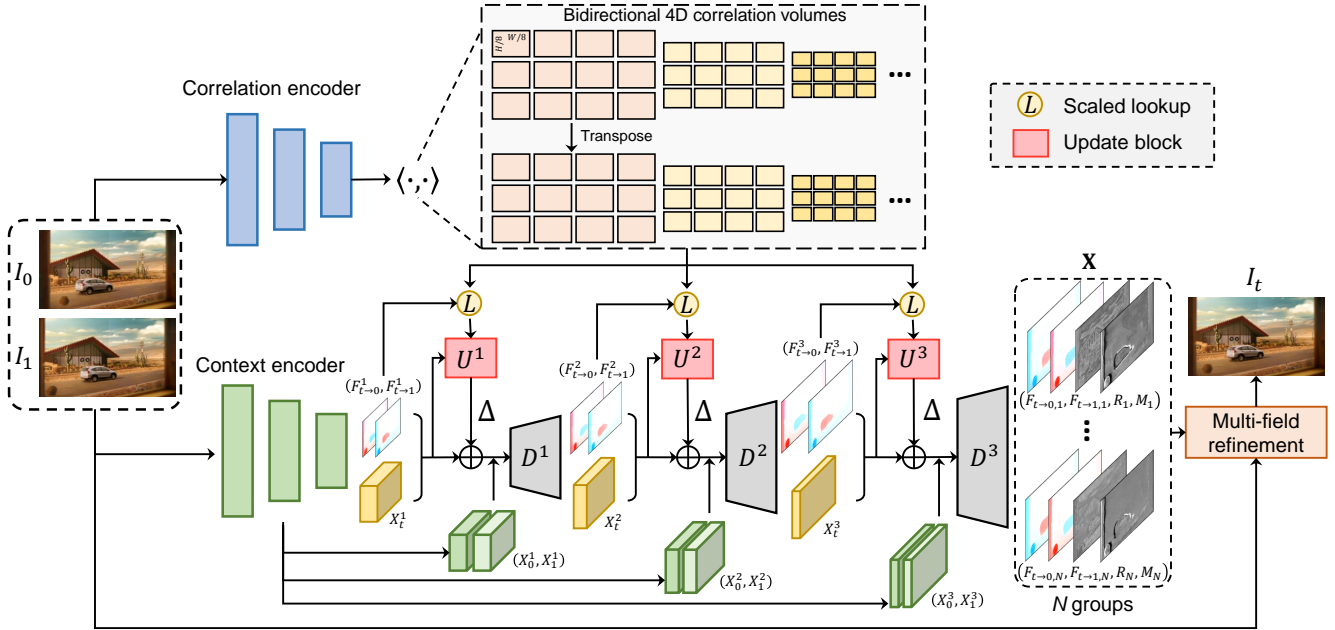


Figure 3. Architecture overview of the proposed AMT. Firstly, the input frames are sent to the correlation encoder to extract features, which are used to construct bidirectional correlation volumes. Then, the context encoder extracts pyramid features of visible frames and generates initial bilateral flows and interpolated intermediate feature. Next, we use bilateral flows to retrieve bidirectional correlations for jointly updating flow fields and the intermediate feature at each level. Finally, we generate multiple groups of flow fields, occlusion masks, and residuals based on the coarse estimate for interpolating the intermediate frame.

ability in motion modeling, which guarantees the consistency of flows on the coarse scale. At the finest scale, we employ multi-field refinement to ensure the diversity for the flow regions that need to be task-specific.

Cost Volume: Cost volume is introduced as a representation of matching costs in numerous vision tasks [13, 21, 26]. In the deep learning era, the concept of cost volume is also proved to be effective in optical flow estimation [17, 20, 60, 61, 75]. Among these works, the most influential ones are PWC-Net [60] and RAFT [61]. In VFI, the existing methods [22, 46, 47, 69] attempt to introduce the cost volume following the scheme of PWC-Net. However, those methods not only search the cost volume in a local region but also depend on inaccurate features warped from reference ones, resulting in a limited performance gain from the cost volume. Instead, the proposed AMT is based on RAFT, which enlarges the search space by iteratively updating the flow field with all-pairs correlation, and only constructs cost volumes between the visible frames. Besides, we involve many *novel* and *task-specific* designs beyond RAFT. The details are described in Sec. 3 and our supplement.

3. Method

Given a pair of input frames (I_0, I_1) , we aim to synthesize an intermediate frame I_t at a target time step t , where $0 < t < 1$. Our AMT is a one-stage flow-based method, in which bilateral flows and the interpolated intermediate

feature are updated and upsampled jointly. As shown in Fig. 3, it is composed of three main components: 1) an encoder for extracting features and initial bilateral flows simultaneously, 2) multi-scale bidirectional correlation volumes for jointly updating bilateral flows and intermediate features at coarse scales, and 3) a multi-field refinement operator for interpolating the target frame with multiple flow groups at the finest scale. Benefiting from such designs, the estimated motion vectors at coarser scales are generally consistent with the ground truth displacements. Meanwhile, they are diverse in fine-grained details at the finest scale, which meets the requirement of *task-oriented* flow. These designs also enable our AMT to capture large motions and successfully handle occlusion regions with high efficiency.

3.1. Initial Flow and Feature Extraction

We employ two separate feature extractors. They are applied to the input pair (I_0, I_1) , but for different purposes. The first is the *correlation encoder*, which maps the input frames to a pair of dense features for constructing bidirectional correlation volumes. We can obtain the pair of features $\mathbf{g}_0, \mathbf{g}_1 \in \mathbb{R}^{H/8 \times W/8 \times D}$ at 1/8 the input image resolution with D channels.

The second is the *context encoder*, which outputs the initial interpolated intermediate feature X_t^1 and predicts the initial bilateral flows $F_{t \rightarrow 0}^1$ and $F_{t \rightarrow 1}^1$. Their spatial resolution is the same as the output of the correlation encoder. Besides, the pyramid features $\{X_0^l, X_1^l \mid l \in \{1, 2, 3\}\}$ for

frames I_0, I_1 are extracted by context encoder for further progressive warping. The architectural details of them can be found in our supplement.

3.2. All-Pairs Correlation

Bidirectional Correlation Volumes: Similar to RAFT [61], we compute the dot-product similarities between all pairs of features vectors for constructing a 4D correlation volume. Given the pair of features $\mathbf{g}_0, \mathbf{g}_1$, we can obtain the correlation volume \mathbf{C} through:

$$\mathbf{C}_{ijkl} = \sum_h \mathbf{g}_{0,ijh} \cdot \mathbf{g}_{1,klh}, \quad \mathbf{C} \in \mathbb{R}^{\frac{H}{8} \times \frac{W}{8} \times \frac{H}{8} \times \frac{W}{8}} \quad (1)$$

For further measuring similarities across scales, the last two dimensions of the correlation volume are downsampled by a repeated 2D average pooling layer with a kernel size of 2 and a stride of 2. We thus obtain a 4-level correlation pyramid $\{\mathbf{C}_1, \mathbf{C}_2, \mathbf{C}_3, \mathbf{C}_4\}$.

However, the correlation pyramid in RAFT is *unidirectional*. It only reflects multi-scale correspondences from I_0 to I_1 . We thereby term it as the forward correlation pyramid. The unidirectional correspondence is insufficient for the VFI task, as the motions are usually asymmetric [47, 72]. Instead of recomputing the matrix multiplication, we directly transpose the correlation volume \mathbf{C} to represent the correspondence in the opposite direction. After obtaining the transposed correlation volume \mathbf{C}^T , we perform the same pooling operation to form the backward correlation pyramid $\{\mathbf{C}_1^T, \mathbf{C}_2^T, \mathbf{C}_3^T, \mathbf{C}_4^T\}$. Note that the bidirectional correlation volumes only need to compute once. The compact global representations assist our network in being aware of large motions in an efficient way.

Correlation Scaled Lookup: After constructing the bidirectional correlation volumes, we intend to query correlation feature maps using estimated bilateral flows $F_{t \rightarrow 0}^l$ and $F_{t \rightarrow 1}^l$. In RAFT, the lookup operation can be directly performed since its estimated flow and the correlation volume share an *identical* coordinate system. For example, the motion $F_{0 \rightarrow 1}$ from frame 0 to frame 1 and the corresponding correlation volume are all based on the coordinate system of the frame 0. Thus, the correlation feature maps can be correctly sampled by the matched flow field. However, for frame interpolation, we can only build correlation volumes from visible reference frames (*i.e.*, I_0, I_1) but estimate the flows (*i.e.*, $F_{t \rightarrow 0}^l$ and $F_{t \rightarrow 1}^l$) of an invisible intermediate frame I_t . So there exists a *mismatch between coordinate systems*, which causes unfaithful correlation lookups and further influences the updating of the flows. A straightforward solution to this problem is transferring bilateral flows $F_{t \rightarrow 0}^l$ and $F_{t \rightarrow 1}^l$ to bidirectional flows $F_{0 \rightarrow 1}^l$ and $F_{1 \rightarrow 0}^l$.

To achieve this goal, we simply scale the estimated bilateral flows based on locally smooth motion assumption [23, 46, 47]. Specifically, we assume the moving objects

are partially overlap within a small time interval. Thus, the bilateral flows and bidirectional flows at the same position are generally consistent in direction but different in magnitude. So the bidirectional flows $F_{0 \rightarrow 1}^l$ and $F_{1 \rightarrow 0}^l$ can be approximated by:

$$F_{0 \rightarrow 1}^l = \frac{1}{1-t} F_{t \rightarrow 1}^l, \quad F_{1 \rightarrow 0}^l = \frac{1}{t} F_{t \rightarrow 0}^l. \quad (2)$$

Subsequently, a lookup operation analogous to that in RAFT performs on bidirectional correlation volumes through approximated bidirectional flows. We construct two lookup windows centered by bidirectional flows with a predefined radius. The lookup operations in the windows are conducted on all levels of the bidirectional correlation pyramids. The retrieved bidirectional correlations are concatenated into one features map for further jointly updating bilateral flows and the interpolated intermediate feature.

Updating with Retrieved Correlations: While RAFT updates and maintains the flow prediction at a single resolution, we predict the bilateral flows in a coarse-to-fine manner following most flow-based VFI methods [18, 28, 41, 42]. This is because that the features of the input pair need to be progressively warped based on the latest flow predictions for generating a faithful intermediate feature. Given the reciprocal relationship between bilateral flow fields and intermediate features in VFI task [18, 28, 35], we also update and upsample the intermediate feature along with the intermediate motions.

Specifically, during the update stage at each spatial level l , we employ an update block to jointly predict the residuals of the bilateral flow fields $F_{t \rightarrow 0}^l, F_{t \rightarrow 1}^l$ and the interpolated intermediate feature X_t^l based on the retrieved bidirectional correlations. In each update block, the bidirectional correlation features and bilateral flows are first passed through two convolutional layers. Then, they are concatenated with the interpolated intermediate feature and injected into two convolutional layers instead of a cumbersome GRU unit in RAFT. Finally, the output features are sent to two separate heads for predicting bilateral flow residuals $\Delta F_{t \rightarrow 0}^l, \Delta F_{t \rightarrow 1}^l$ and an interpolated feature residual ΔX_t^l . Each head is formed by two convolutional layers.

Note that the spatial dimension of the retrieved correlation features is the same as the first two dimensions of the correlation volume (*i.e.*, $\frac{H}{8} \times \frac{W}{8}$) but is different from that of the intermediate features and motions on upper levels. We thus need to downscale the flow fields and the intermediate feature accordingly before feeding them into the update block and upsample the predicted residuals for updating. Through downscaling, the update block works at a low-resolution space, leading to promising efficiency. The updated intermediate feature \hat{X}_t^l can be formulated as: $\hat{X}_t^l = X_t^l + \Delta X_t^l$, where ΔX_t^l is the output content residual of the update block. The updated bilateral flows $\hat{F}_{t \rightarrow 0}^l, \hat{F}_{t \rightarrow 1}^l$ can be obtained following the same rule.

We employ the updated bilateral flows to warp the features X_0^l, X_1^l of the input frames. Let \hat{X}_0^l, \hat{X}_1^l denote the warped features. The warped features, the updated bilateral flows, and the updated intermediate feature are concatenated together and then fed into the l -th decoder. The l -th decoder D^l predicts the upsampled bilateral flows $F_{t \rightarrow 0}^{l+1}, F_{t \rightarrow 1}^{l+1}$ and the intermediate feature X_t^{l+1} simultaneously as follows:

$$[F_{t \rightarrow 0}^{l+1}, F_{t \rightarrow 1}^{l+1}, X_t^{l+1}] = D^l([\hat{X}_0^l, \hat{X}_1^l, \hat{F}_{t \rightarrow 0}^l, \hat{F}_{t \rightarrow 1}^l, \hat{X}_t^l]). \quad (3)$$

Specially, the Eqn. (3) does not consider the last decoder D^3 , which is responsible for generating multiple flow fields and occlusion masks for task-specific usage. The architecture details of each decoder are listed in our supplement.

3.3. Multi-Field Refinement

In flow-based VFI methods, the common formulation for interpolating the final intermediate frame is:

$$I_t = M \odot \mathcal{W}(I_0, F_{t \rightarrow 0}) + (1 - M) \odot \mathcal{W}(I_1, F_{t \rightarrow 1}) + R, \quad (4)$$

where \mathcal{W} denotes the backward warping operation, \odot means the element-wise multiplication. M is an estimated occlusion mask which ranges from 0 to 1. $F_{t \rightarrow 0}$ and $F_{t \rightarrow 1}$ are final predictions of bilateral flows. R is the estimated residual. Such formulation considers temporal consistency and occlusion reasoning, synthesizing the intermediate frame efficiently. However, only predicting one pair of flow fields ignores that each location in the occlusion areas has many potential pixel candidates, restricting the solution set for interpolation in a tight space.

Based on previously predicted coarse flows, which are generally consistent with the ground truth displacements, we derive multiple fine-grained flow fields for task-specific usage. We also jointly estimate a residual content and an occlusion mask for each pair of optical flow. This process can be formulated as:

$$\mathbf{X} = D^3([\hat{X}_0^3, \hat{X}_1^3, \hat{F}_{t \rightarrow 0}^3, \hat{F}_{t \rightarrow 1}^3, \hat{X}_t^3]), \quad (5)$$

$$\mathbf{X} = \{F_{t \rightarrow 0, n}, F_{t \rightarrow 1, n}, M_n, R_n | n \in \{1, 2, \dots, N\}\},$$

where N denotes the total number of output groups. $(F_{t \rightarrow 0, n}, F_{t \rightarrow 1, n})$, M_n , and R_n are the n -th estimated bilateral flows, occlusion mask, and residual content, respectively. Notably, Eqn. (5) can be easily implemented by enlarging the output channels of the last decoder according to the number of flow pairs, which ensures efficiency. The final intermediate frame can be obtained by:

$$I_t = \mathcal{C}([I_t^1, \dots, I_t^N]), \quad (6)$$

where the n -th interpolated frame I_t^n can be obtained by Eqn. (4) with corresponding output group. We stack two convolutional layers (denoted as \mathcal{C}) for adaptively merging candidate frames and refining the final results. The analyses of multiple flow fields are detailed in Sec. 4.4.2.

3.4. Loss Functions

There are three losses involved in our AMT. To better predict task-oriented flows, we employ flow distillation loss \mathcal{L}_{flow} in IFRNet [28], which concentrates more on the flow regions that are easy to be reconstructed, but slightly penalizes the regions that are difficult to recover. This loss is applied on updated multi-scale flow fields except for the finest flow predictions left for fully task-specific usage. The Charbonnier loss [4] \mathcal{L}_{char} and the census loss [38] \mathcal{L}_{css} are used to supervise the content generation of the interpolated frame. The former measures the pixel-wise errors between the ground truth intermediate frame I_t^{GT} and the generated one I_t , and the latter calculates the soft Hamming distance between census-transformed image patches of I_t^{GT} and I_t .

The full objective can be defined as:

$$\mathcal{L} = \lambda_{char} \mathcal{L}_{char} + \lambda_{css} \mathcal{L}_{css} + \lambda_{flow} \mathcal{L}_{flow}, \quad (7)$$

where λ_{char} , λ_{css} , and λ_{flow} are weights for each loss.

4. Experiments

4.1. Training Details

We train AMT on Vimeo90K [73] training set for 300 epochs with AdamW [36] optimizer on 2 NVIDIA RTX 3090 GPUs. The total batch size is 24, and the learning rate decay follows the cosine attenuation schedule from 2×10^{-4} to 2×10^{-5} . We follow the augmentation pipeline including random flipping, rotating, reversing sequence order, and random cropping patches with size 224×224 in IFRNet [28]. The flow predictions from the pre-trained LiteFlowNet [19] are served as the pseudo ground truth label for supervising the intermediate flows. λ_{char} , λ_{css} , and λ_{flow} are set as 1, 1, and 0.002, respectively. The code implemented by MindSpore framework is also provided.

4.2. Benchmarks

We evaluate our AMT on various benchmarks containing diverse motion scenes for a comprehensive comparison. PSNR and SSIM [66] as common evaluation metrics are utilized for comparison. The statistics of benchmarks used in the main paper are presented as follows.

Vimeo90K [73]: Vimeo90K is the most commonly used evaluation benchmark in recent VFI literature. There are 3,782 triplets of 448×256 resolution in the test part.

UCF101 [58]: UCF101 dataset contains videos with various human actions, and we adopt the test partition in DVF [35], which consists of 379 triplets of 256×256 size.

SNU-FILM [8]: SNU-FILM dataset contains 1,240 frame triplets, whose width ranges from 368 to 720 and height ranges from 384 to 1280. With respect to motion magnitude, it is partitioned into four exclusive parts, namely Easy, Medium, Hard, and Extreme.

Xiph [39]: Xiph dataset, consisting of eight video clips with a 4K resolution, was originally proposed by Niklaus

Method	Vimeo90K [73]	UCF101 [58]	SNU-FILM [8]				Xiph [39]		Latency (ms/f)	Params (M)	FLOPs (T)
			Easy	Medium	Hard	Extreme	2K	4K			
AdaCoF [29]	34.38/0.972	35.20/0.970	39.85/0.991	35.08/0.976	29.47/0.925	24.31/0.844	34.86/0.928	31.68/0.870	52	21.8	0.36
M2M-VFI [15]	35.49/0.978	<u>35.32/0.970</u>	39.66/0.991	35.74/0.980	30.32/0.936	25.07/0.860	36.44/0.943	33.92/0.899	40	7.6	0.26
RIFE [18]	35.65/0.978	<u>35.28/0.969</u>	40.06/0.991	35.75/0.979	30.10/0.933	24.84/0.853	36.19/0.938	33.76/0.894	29	9.8	0.20
IFRNet-S [28]	35.59/0.979	35.28/0.969	<u>39.96/0.991</u>	35.92/0.979	30.36/0.936	25.05/0.858	35.87/0.936	33.80/0.891	25	2.8	0.12
IFRNet-B [28]	<u>35.80/0.979</u>	35.29/0.969	40.03/0.991	<u>35.94/0.979</u>	<u>30.41/0.936</u>	<u>25.05/0.859</u>	36.00/0.936	<u>33.99/0.893</u>	30	5.0	0.21
AMT-S	35.97/0.983	35.35/0.971	39.95/0.994	35.98/0.983	30.60/0.940	25.30/0.865	<u>36.11/0.940</u>	34.29/0.901	51	3.0	0.12
ToFlow [73]	33.73/0.968	34.58/0.967	39.08/0.989	34.39/0.974	28.44/0.918	23.39/0.831	33.93/0.922	30.74/0.856	88	1.4	0.62
DAIN [1]	34.71/0.976	34.99/0.968	39.73/0.990	35.46/0.978	30.17/0.934	25.09/0.858	35.95/0.940	33.49/0.895	664	24.0	5.51
CAIN [8]	34.78/0.974	35.00/0.969	39.95/0.990	35.66/0.978	29.93/0.930	24.80/0.851	35.21/0.937	32.56/0.901	71	42.8	1.29
BMBC [46]	35.01/0.976	35.15/0.969	39.90/0.990	35.31/0.977	29.33/0.927	23.92/0.843	32.82/0.928	31.19/0.880	2234	11.0	2.50
ABME [47]	<u>36.22/0.981</u>	<u>35.41/0.970</u>	39.59/0.990	35.77/0.979	30.58/0.937	25.42/0.864	36.53/0.944	33.73/0.901	560	18.1	1.30
IFRNet-L [28]	36.20/0.981	<u>35.42/0.970</u>	40.10/0.991	36.12/0.980	30.63/0.937	25.27/0.861	36.21/0.937	<u>34.25/0.895</u>	80	19.7	0.79
AMT-L	36.35/0.982	35.42/0.970	39.95/0.991	<u>36.09/0.981</u>	30.75/0.938	<u>25.41/0.864</u>	<u>36.27/0.940</u>	34.49/0.903	116	12.9	0.58
VFIFormer [37]	<u>36.50/0.982</u>	<u>35.43/0.970</u>	40.13/0.991	36.09/0.980	30.67/0.938	25.43/0.864	OOM	OOM	1293	24.1	47.71
EMA-VFI [†] [76]	36.50/0.980	35.42/0.970	39.58/0.989	35.86/0.979	30.80/0.938	25.59/0.864	36.74/0.944	<u>34.55/0.906</u>	211	66.0	0.91
AMT-G	36.53/0.982	35.45/0.970	<u>39.88/0.991</u>	36.12/0.981	<u>30.78/0.939</u>	<u>25.43/0.865</u>	<u>36.38/0.941</u>	34.63/0.904	250	30.6	2.07

Table 1. Quantitative comparison with SOTA methods. We divide the existing methods into three groups, according to the computational complexity. For each group, the best result is shown in **bold**, and the second best is underlined. “OOM” denotes the out-of-memory issue when evaluating on an NVIDIA RTX 3090 GPU. [†] means we disable the test-time augmentation [18] for a fair comparison.

et al. [42]. Following their original evaluation setting, we reform this dataset to include “2K” version, obtained by downscaling original frames, and “4K” version, created by center-cropping 2K patches.

Except for these datasets, we provide the comparisons of multi-frame interpolation in the supplement.

4.3. Comparison with the SOTAs

We compare our AMT with the state-of-the-art (SOTA) methods, including ToFlow [73], DAIN [1], CAIN [8], AdaCoF [29], BMBC [46], RIFE [18], ABME [47], M2M-VFI [15], IFRNet [28], VFIFormer [37], and EMA-VFI [76]. We utilize the code provided by IFRNet [28] for benchmarks. The inference latency is the average running time of a method on 1280×720 resolution for 1000 iterations on an NVIDIA RTX 3090 GPU. To ensure a fair comparison, we group the SOTA methods into three categories based on their theoretical computational complexity. We then develop three models, called AMT-S, AMT-L, and AMT-G, for each group.

Quantitative Comparison. As shown in Tab. 1, our small model AMT-S achieves the best results among efficient VFI methods on almost all benchmarks, especially for challenging settings. Specifically, Our AMT-S outperforms the previous state-of-the-art method in effective VFI, IFRNet-B [28], by 0.17dB on Vimeo90K while using only about 60% of its parameters and FLOPs. This gap becomes more obvious on the Hard and Extreme partitions in SNU-FILM, revealing the strong ability of our AMT in modeling large motions. For the large scale setting, our AMT-L shows highly competitive results in contrast to the previous SOTA method IFRNet-L [28], with about 65% parameters and 75% FLOPs of it. In terms of inference speed, our

method is comparable to IFRNet. Besides, our convolution-based model competes favorably compared to the SOTA Transformer-based models (*i.e.*, VFIFormer [37] and EMA-VFI [76]) in terms of accuracy and efficiency. Specifically, our AMT-G outperforms them in most cases, particularly when evaluated using the SSIM metric. Notably, our model achieves about $5\times$ faster inference speed than VFIFormer and has only half the number of parameters of EMA-VFI. It is important to note that VFIFormer requires a two-stage training pipeline and 600 training epochs, while our model only requires 300 epochs. Additionally, EMA-VFI introduces a warm-up technique during training, which our method does not utilize. We observe that the performance of our method is saturated except for the Vimeo90K dataset after increasing the scale of the model to a huge version, which may indicate the overfitting problem.

Qualitative Comparison. In Fig. 4, we select the representative hallucination-based, kernel-based, and flow-based methods, including CAIN [8], AdaCoF [29], ABME [47], RIFE [18], and IFRNet(-B/-L) [28]. We compare them with our AMT on SNU-FILM [8] (Hard) dataset for visual comparison. It can be seen that previous VFI methods fail to provide sharp edges of moving objects, especially when the motion is complex. Due to our thorough consideration of VFI-oriented flows, our AMT synthesizes the content at motion boundaries faithfully and generates plausible textures with fewer artifacts. When the background objects are heavily occluded by the foreground unilaterally, our AMT can still obtain guidance from the reference frame in another direction, while other methods are unable to synthesize these occluded objects. We provide more comparisons in the supplement.

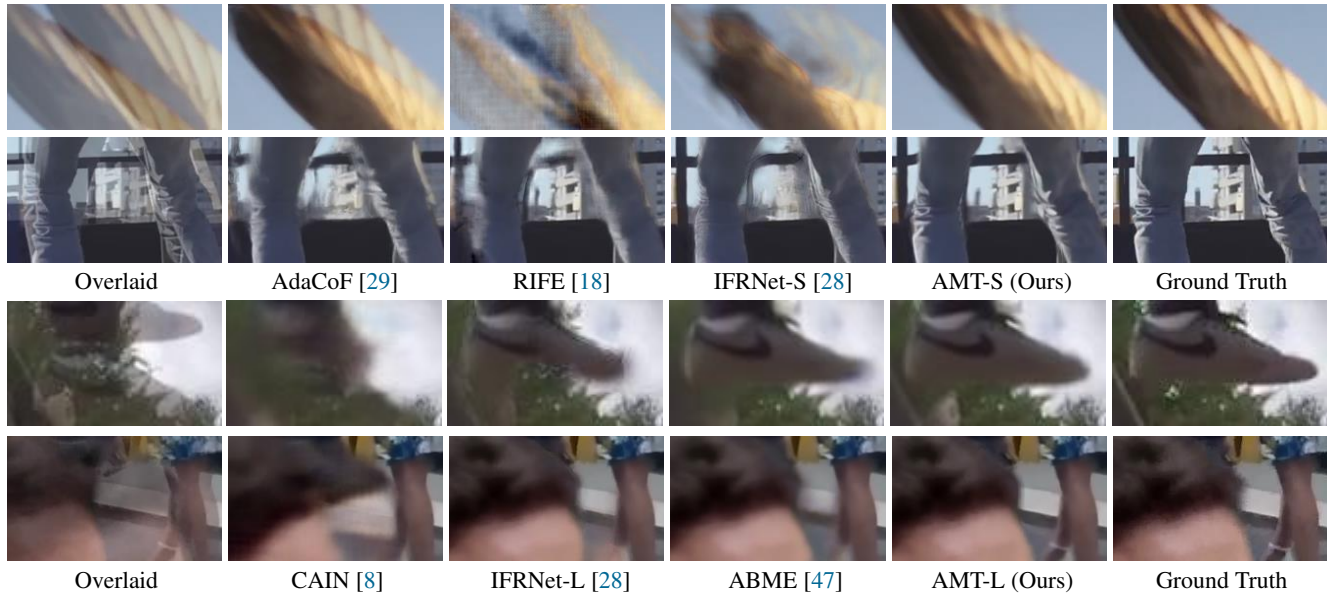


Figure 4. Qualitative results from different VFI methods. We divide these methods into two groups by computational cost. Our AMT-S and AMT-B synthesize precise boundaries of the objects with large motion and can reconstruct occluded regions with high fidelity.

Case	Vimeo	Hard	Extreme
w/o Corr. Enc.	35.76	30.49	25.22
Unidir. CV	35.93	30.34	25.18
PWC CV	35.61	30.48	25.16
Full Model	35.97	30.60	25.30

(a) **Correlation volume (CV) design.** We remove the correlation encoder ('w/o Corr. Enc.'), build a unidirectional CV ('Unidir. CV'), and build PWC-like [51] CV ('PWC CV') for ablations, respectively.

1st	2nd	3rd	Vimeo	Hard	Extreme
			35.60	30.39	25.06
✓			35.84	30.55	25.19
✓	✓		35.92	30.58	25.28
✓	✓	✓	35.97	30.60	25.30
		single-scale	35.95	30.50	25.22

(d) **Cross-scale update.** We investigate the impact of the update at different levels.

Table 2. Ablation experiments of AMT on Vimeo90K [73] and SNU-FILM [8] (Hard, Extreme) dataset. We report the PSNR values of these variants, and the best result is shown in **bold**. The default setting is marked in **gray**.

4.4. Ablation Study

We conduct ablations to verify the effectiveness of two key components (*i.e.*, all-pairs correlation and multi-field refinement) in our AMT. All ablated versions are based on the AMT-S and evaluated on Vimeo90K [73] and the Hard and Extreme partitions of SNU-FILM [8].

4.4.1 All-Pairs Correlation

Volume Designs. As illustrated in Tab. 2a, our bidirectional correlation volumes thoroughly consider the correspondences between input frames for the VFI task, leading to a better performance than the unidirectional one. Besides, using an exclusive encoder (*i.e.*, correlation encoder)

Lookup	Init	Vimeo	Hard	Extreme
Initial meshgrid		35.92	30.52	25.23
RAFT Flow		35.93	30.34	25.18
Scaled Zero		35.97	30.56	25.26
Scaled Flow		35.97	30.60	25.30

(b) **Lookup strategy.** We investigate the initial meshgrid, RAFT-like [61] lookup ('RAFT'), and the proposed lookup ('Scaled') variants. We also investigate whether we use bilateral flows to perform an initial lookup.

No.	Vimeo	Hard	Extreme	FLOPs (G)
1	35.84	30.52	25.25	116
3	35.97	30.60	25.30	121
5	36.00	30.63	25.33	127
7	36.01	30.57	25.25	135

(e) **Number of fields.** We investigate different numbers of flow pairs.

Case	Vimeo	Hard	Extreme
Vanilla Guide	35.95	30.53	25.21
w/o Update	35.96	30.52	25.22
Full Model	35.97	30.60	25.30

(c) **Content update.** We investigate the content update by using features from visible frames as guidance ('Vanilla Guide') and discarding the content update ('w/o Update'), respectively.

Case	Vimeo	Hard	Extreme
w/o Residual	35.87	30.57	25.27
w/o Refine	35.89	30.51	25.19
Full Model	35.97	30.60	25.30

(f) **Multi-field combination.** We investigate the residual component in Eqn. (4) and the refinement step in Eqn. (6).

for building the correlation volumes is necessary. We can observe that the performance heavily drops when we utilize features from the context encoder to construct the correlation volumes. We also try to build the correlation volumes following PWC-Net [60]. This variant performs worse than any other one, for its partial correlation volume limits the ability in modeling motion information sufficiently.

Lookup Strategy. As shown in Tab. 2b, we can observe an obvious performance drop while utilizing the vanilla lookup strategy in RAFT [61]. For large motions, its performance is even worse than the one that directly uses the initial meshgrid, which indicates this strategy provides unfaithful correlation information for flow updates. After we project the

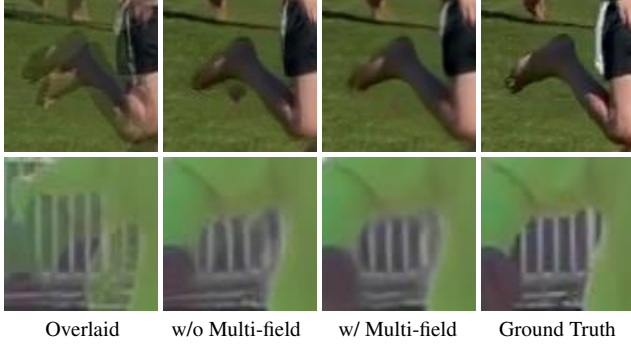


Figure 5. Effect of multi-field refinement. Multi-field refinement helps the network recover occluded regions better.

flows by scaling, the correlation volumes and flows share an identical coordinate system, and the network takes advantage of the correct lookup process. Besides, the initial flow pair from the context encoder gives a good initial point for further lookup, which brings a performance gain.

Content Update. In our AMT, each update block receives the intermediate content features as the context guidance and updates it along with bilateral flows. If we replace the context guidance with features from visible frames, the ambiguous information will be introduced, leading to a performance drop, as shown in Tab. 2c. Besides, we only keep one head in each update block for only updating the flow fields without updating the intermediate feature, resulting in the decrease of PSNR values on large motions. It demonstrates that all-pairs correlation is not only helpful for updating flows but also for updating content.

Update Strategy. As shown in Tab. 2d, all updates across levels are effective in our cross-scale update strategy. It is worth noticing that if we discard all updates, which is equivalent to a model without all-pairs correlation, the PSNR value will decrease dramatically. This demonstrates *the effectiveness of all-pairs correlation* in our AMT. Besides, only updating on the 1-st scale with $3\times$ iterations degrades the performance. The fact indicates that the cross-scale update strategy can take full advantage of progressively refined content features, leading to better motion modeling.

4.4.2 Multi-field Refinement

Number of Flow Fields. Tab. 2e illustrates the performance gain with respect to the number of flow fields. We observe that just using three pairs of flows bring a notable performance gain, which reveals that ensuring the diversity of flow fields is significant for VFI-oriented usage. The PSNR values rises in pace with the increase of field number until 7 pairs, which indicates saturation. We employ 3 pairs in our small model for efficiency (*i.e.*, AMT-S) and use 5 pairs in the larger models for better performance. In Fig. 5, we investigate the effect of multi-field refinement on occlusion handling. The results indicate that after employing multi-field refinement, our AMT can synthesize the background

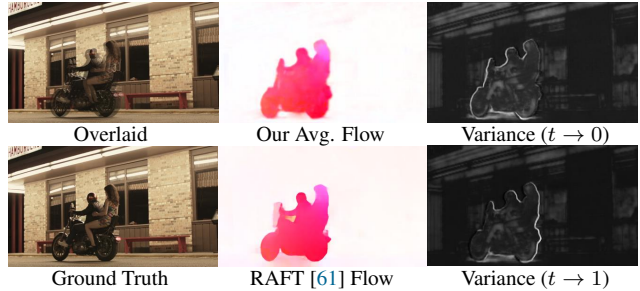


Figure 6. Visualizations of average and variance map of three flow pairs. We provide RAFT [61] flow for reference.

occluded by the foreground with more consistent textures.

Multi-Field Combination. We investigate a variant that removes the residual component for each candidate frame in Eqn. (4) but estimates the residual part in the final interpolation result. As shown in Tab. 2f, the results of this variant underperform the original setting, which indicates we need to compensate details for each frame candidate separately. Besides, if we replace the convolution operators in Eqn. (6) with an average operation, the performance will be degraded (see Tab. 2f). This indicates that it is important for our AMT to perform an adaptive fusion and refinement.

Discussion. For further discussion, we visualize the mean and deviation of three estimated flow pairs. The results are shown in Fig. 6. On the one hand, our average flow is generally consistent with the flow estimated from RAFT [61], which approximates to the ground truth displacements. On the other hand, we observe that the major diversities of flows are at the motion boundaries and in the regions with rich textures. This indicates that these regions need to involve more potential pixel candidates for reconstruction. Through these visualizations, we see that our method generate promising task-oriented flows, *generally consistent with the ground truth optical flows but diverse in local details.*

5. Conclusion

Following the property of task-oriented flow, we have introduced All-pairs Multi-field Transforms (AMT) for efficient frame interpolation. It contains two essential designs, including all-pairs correlation and multi-field refinement. Through the two designs, our method could effectively handle large motions and occluded regions during frame interpolation and achieve state-of-the-art performance on various benchmarks with high efficiency.

Acknowledgment: This work is funded by the NSFC (NO. 62176130), Fundamental Research Funds for the Central Universities (Nankai University, NO. 63223050), China Postdoctoral Science Foundation (NO.2021M701780). This work is partially supported by Ascend AI Computing Platform and CANN (Compute Architecture for Neural Networks). We are also sponsored by CAAI-Huawei MindSpore Open Fund.

References

- [1] Wenbo Bao, Wei-Sheng Lai, Chao Ma, Xiaoyun Zhang, Zhiyong Gao, and Ming-Hsuan Yang. Depth-aware video frame interpolation. In *CVPR*, 2019. 2, 6, 11
- [2] Wenbo Bao, Wei-Sheng Lai, Xiaoyun Zhang, Zhiyong Gao, and Ming-Hsuan Yang. Memc-net: Motion estimation and motion compensation driven neural network for video interpolation and enhancement. *IEEE TPAMI*, 2018. 1, 2
- [3] Kelvin C.K. Chan, Shangchen Zhou, Xiangyu Xu, and Chen Change Loy. BasicVSR++: Improving video super-resolution with enhanced propagation and alignment. In *CVPR*, 2022. 2
- [4] Pierre Charbonnier, Laure Blanc-Feraud, Gilles Aubert, and Michel Barlaud. Two deterministic half-quadratic regularization algorithms for computed imaging. In *ICIP*, 1994. 5
- [5] Xianhang Cheng and Zhenzhong Chen. Video frame interpolation via deformable separable convolution. In *AAAI*, 2020. 2
- [6] Xianhang Cheng and Zhenzhong Chen. Multiple video frame interpolation via enhanced deformable separable convolution. *IEEE TPAMI*, 2021. 2
- [7] Zhixiang Chi, Rasoul Mohammadi Nasiri, Zheng Liu, Yuanhao Yu, Juwei Lu, Jin Tang, and Konstantinos N Plataniotis. Error-aware spatial ensembles for video frame interpolation. *arXiv preprint arXiv:2207.12305*, 2022. 2
- [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*, 2020. 2, 5, 6, 7, 12, 13, 18, 19, 20, 21, 22
- [9] Jifeng Dai, Haozhi Qi, Yuwen Xiong, Yi Li, Guodong Zhang, Han Hu, and Yichen Wei. Deformable convolutional networks. In *CVPR*, 2017. 2
- [10] Tianyu Ding, Luming Liang, Zhihui Zhu, and Ilya Zharkov. Cdfi: Compression-driven network design for frame interpolation. In *CVPR*, 2021. 2
- [11] John Flynn, Ivan Neulander, James Philbin, and Noah Snavely. Deepstereo: Learning to predict new views from the world’s imagery. In *CVPR*, 2016. 1
- [12] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 11
- [13] Asmaa Hosni, Christoph Rhemann, Michael Bleyer, Carsten Rother, and Margrit Gelautz. Fast cost-volume filtering for visual correspondence and beyond. *IEEE TPAMI*, 2012. 3
- [14] Jie Hu, Li Shen, and Gang Sun. Squeeze-and-excitation networks. In *CVPR*, 2018. 2
- [15] Ping Hu, Simon Niklaus, Stan Sclaroff, and Kate Saenko. Many-to-many splatting for efficient video frame interpolation. In *CVPR*, 2022. 1, 6, 12, 13
- [16] Xiaotao Hu, Zhewei Huang, Ailin Huang, Jun Xu, and Shuchang Zhou. A dynamic multi-scale voxel flow network for video prediction. In *CVPR*, 2023. 2
- [17] Zhaoyang Huang, Xiaoyu Shi, Chao Zhang, Qiang Wang, Ka Chun Cheung, Hongwei Qin, Jifeng Dai, and Hongsheng Li. FlowFormer: A transformer architecture for optical flow. In *ECCV*, 2022. 3, 11
- [18] Zhewei Huang, Tianyuan Zhang, Wen Heng, Boxin Shi, and Shuchang Zhou. Real-time intermediate flow estimation for video frame interpolation. In *ECCV*, 2022. 1, 2, 4, 6, 7, 11, 13, 17, 19, 21
- [19] Tak-Wai Hui, Xiaoou Tang, and Chen Change Loy. Lite-FlowNet: A Lightweight Convolutional Neural Network for Optical Flow Estimation. In *CVPR*, 2018. 5
- [20] Junhwa Hur and Stefan Roth. Iterative residual refinement for joint optical flow and occlusion estimation. In *CVPR*, 2019. 3
- [21] Sangryul Jeon, Seungryong Kim, Dongbo Min, and Kwanghoon Sohn. Pam: Pyramidal affine regression networks for dense semantic correspondence. In *ECCV*, 2018. 3
- [22] Zhaoyang Jia, Yan Lu, and Houqiang Li. Neighbor correspondence matching for flow-based video frame synthesis. In *ACM MM*, 2022. 3
- [23] Huaizu Jiang, Deqing Sun, Varun Jampani, Ming-Hsuan Yang, Erik Learned-Miller, and Jan Kautz. Super slo-mo: High quality estimation of multiple intermediate frames for video interpolation. In *CVPR*, 2018. 1, 2, 4, 11
- [24] Shihao Jiang, Dylan Campbell, Yao Lu, Hongdong Li, and Richard Hartley. Learning to estimate hidden motions with global motion aggregation. In *ICCV*, 2021. 11
- [25] Tarun Kalluri, Deepak Pathak, Manmohan Chandraker, and Du Tran. Flavr: Flow-agnostic video representations for fast frame interpolation. *arXiv preprint arXiv:2012.08512*, 2020. 2
- [26] Alex Kendall, Hayk Martirosyan, Saumitro Dasgupta, Peter Henry, Ryan Kennedy, Abraham Bachrach, and Adam Bry. End-to-end learning of geometry and context for deep stereo regression. In *ICCV*, 2017. 3
- [27] Soo Ye Kim, Jihyong Oh, and Munchurl Kim. Fisr: Deep joint frame interpolation and super-resolution with a multi-scale temporal loss. In *AAAI*, 2020. 2
- [28] Lingtong Kong, Boyuan Jiang, Donghao Luo, Wenqing Chu, Xiaoming Huang, Ying Tai, Chengjie Wang, and Jie Yang. Ifrnet: Intermediate feature refine network for efficient frame interpolation. In *CVPR*, 2022. 1, 2, 4, 5, 6, 7, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22
- [29] 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*, 2020. 2, 6, 7, 13, 17, 19, 21
- [30] Zhen Li, Cheng-Ze Lu, Jianhua Qin, Chun-Le Guo, and Ming-Ming Cheng. Towards an end-to-end framework for flow-guided video inpainting. In *CVPR*, 2022. 2
- [31] Zhengqi Li, Simon Niklaus, Noah Snavely, and Oliver Wang. Neural scene flow fields for space-time view synthesis of dynamic scenes. In *CVPR*, 2021. 1
- [32] Jing Lin, Yuanhao Cai, Xiaowan Hu, Haoqian Wang, Youliang Yan, Xueyi Zou, Henghui Ding, Yulun Zhang, Radu Timofte, and Luc Van Gool. Flow-guided sparse transformer for video deblurring. In *ICML*, 2022. 2
- [33] Lahav Lipson, Zachary Teed, and Jia Deng. Raft-stereo: Multilevel recurrent field transforms for stereo matching. In *3DV*, 2021. 11

- [34] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *ICCV*, 2021. 2
- [35] Ziwei Liu, Raymond Yeh, Xiaoou Tang, Yiming Liu, and Aseem Agarwala. Video frame synthesis using deep voxel flow. In *ICCV*, 2017. 1, 2, 4, 5, 11
- [36] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *ICLR*, 2018. 5
- [37] Liying Lu, Ruizheng Wu, Huajia Lin, Jiangbo Lu, and Jiaya Jia. Video frame interpolation with transformer. In *CVPR*, 2022. 1, 2, 6, 13
- [38] Simon Meister, Junhwa Hur, and Stefan Roth. UnFlow: Un-supervised learning of optical flow with a bidirectional census loss. In *AAAI*, 2018. 5
- [39] Christopher Montgomery. Xiph.org video test media (derf’s collection). In *Online*, <https://media.xiph.org/video/derf/>, 1994. 5, 6
- [40] Seungjun Nah, Tae Hyun Kim, and Kyoung Mu Lee. Deep multi-scale convolutional neural network for dynamic scene deblurring. In *CVPR*, 2017. 11
- [41] Simon Niklaus and Feng Liu. Context-aware synthesis for video frame interpolation. In *CVPR*, 2018. 1, 2, 4
- [42] Simon Niklaus and Feng Liu. Softmax splatting for video frame interpolation. In *CVPR*, 2020. 1, 2, 4, 6
- [43] Simon Niklaus, Long Mai, and Feng Liu. Video frame interpolation via adaptive convolution. In *CVPR*, 2017. 2
- [44] Simon Niklaus, Long Mai, and Feng Liu. Video frame interpolation via adaptive separable convolution. In *ICCV*, 2017. 2
- [45] Simon Niklaus, Long Mai, and Oliver Wang. Revisiting adaptive convolutions for video frame interpolation. In *WACV*, 2021. 2
- [46] Junheum Park, Keunsoo Ko, Chul Lee, and Chang-Su Kim. Bmbc: Bilateral motion estimation with bilateral cost volume for video interpolation. In *ECCV*, 2020. 2, 3, 4, 6, 12
- [47] Junheum Park, Chul Lee, and Chang-Su Kim. Asymmetric bilateral motion estimation for video frame interpolation. In *CVPR*, 2021. 1, 3, 4, 6, 7, 12, 13, 18, 20, 22
- [48] Tomer Peleg, Pablo Szekely, Doron Sabo, and Omry Sendik. Im-net for high resolution video frame interpolation. In *CVPR*, 2019. 2
- [49] Feng Liu Qiqi Hou, Abhijay Ghildyal. A perceptual quality metric for video frame interpolation. In *ECCV*, 2022. 1
- [50] Fitsum Reda, Janne Kontkanen, Eric Tabellion, Deqing Sun, Caroline Pantofaru, and Brian Curless. Film: Frame interpolation for large motion. In *ECCV*, 2022. 2, 11
- [51] Zhile Ren, Orazio Gallo, Deqing Sun, Ming-Hsuan Yang, Erik B Sudderth, and Jan Kautz. A fusion approach for multi-frame optical flow estimation. In *WACV*, 2019. 7
- [52] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *MICCAI*, 2015. 1
- [53] Wenzhe Shi, Jose Caballero, Ferenc Huszár, Johannes Totz, Andrew P Aitken, Rob Bishop, Daniel Rueckert, and Zehan Wang. Real-time single image and video super-resolution using an efficient sub-pixel convolutional neural network. In *CVPR*, 2016. 2
- [54] Zhihao Shi, Xiangyu Xu, Xiaohong Liu, Jun Chen, and Ming-Hsuan Yang. Video frame interpolation transformer. In *CVPR*, 2022. 2
- [55] Hyeonjun Sim, Jiyoung Oh, and Munchurl Kim. Xvfi: extreme video frame interpolation. In *ICCV*, 2021. 2
- [56] Uriel Singer, Adam Polyak, Thomas Hayes, Xi Yin, Jie An, Songyang Zhang, Qiyuan Hu, Harry Yang, Oron Ashual, Oran Gafni, et al. Make-a-video: Text-to-video generation without text-video data. *arXiv preprint arXiv:2209.14792*, 2022. 1
- [57] Li Siyao, Shiyu Zhao, Weijiang Yu, Wenxiu Sun, Dimitris Metaxas, Chen Change Loy, and Ziwei Liu. Deep animation video interpolation in the wild. In *CVPR*, 2021. 1
- [58] 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. 5, 6
- [59] Shuo Chen Su, Mauricio Delbracio, Jue Wang, Guillermo Sapiro, Wolfgang Heidrich, and Oliver Wang. Deep video deblurring for hand-held cameras. In *CVPR*, 2017. 11, 12, 16
- [60] Deqing Sun, Xiaodong Yang, Ming-Yu Liu, and Jan Kautz. Pwc-net: Cnns for optical flow using pyramid, warping, and cost volume. In *CVPR*, 2018. 3, 7
- [61] Zachary Teed and Jia Deng. Raft: Recurrent all-pairs field transforms for optical flow. In *ECCV*, 2020. 1, 2, 3, 4, 7, 8, 11, 12
- [62] Zachary Teed and Jia Deng. Raft-3d: Scene flow using rigid-motion embeddings. In *CVPR*, 2021. 11
- [63] Du Tran, Lubomir Bourdev, Rob Fergus, Lorenzo Torresani, and Manohar Paluri. Learning spatiotemporal features with 3d convolutional networks. In *ICCV*, 2015. 2
- [64] Dmitry Ulyanov, Andrea Vedaldi, and Victor Lempitsky. Instance normalization: The missing ingredient for fast stylization. *arXiv preprint arXiv:1607.08022*, 2016. 11
- [65] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, 2017. 2
- [66] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE TIP*, 2004. 5
- [67] Chao-Yuan Wu, Nayan Singhal, and Philipp Krahenbuhl. Video compression through image interpolation. In *ECCV*, 2018. 1
- [68] Xiaoyu Xiang, Yapeng Tian, Yulun Zhang, Yun Fu, Jan P Allebach, and Chenliang Xu. Zooming slow-mo: Fast and accurate one-stage space-time video super-resolution. In *CVPR*, 2020. 1
- [69] Jin Xin, Wu Longhai, Shen Guotao, Chen Youxin, Chen Jie, Koo Jayoon, and Hahm Cheul-hee. Enhanced bi-directional motion estimation for video frame interpolation. *arXiv preprint arXiv:2206.08572*, 2022. 3
- [70] Haofei Xu, Jiaolong Yang, Jianfei Cai, Juyong Zhang, and Xin Tong. High-resolution optical flow from 1d attention and correlation. In *ICCV*, 2021. 11
- [71] Haofei Xu, Jing Zhang, Jianfei Cai, Hamid Reza Tofighi, and Dacheng Tao. Gmflow: Learning optical flow via global matching. In *CVPR*, 2022. 1, 11

- [72] Xiangyu Xu, Li Siyao, Wenxiu Sun, Qian Yin, and Ming-Hsuan Yang. Quadratic video interpolation. In *NeurIPS*, 2019. 2, 4
- [73] Tianfan Xue, Baian Chen, Jiajun Wu, Donglai Wei, and William T Freeman. Video enhancement with task-oriented flow. *IJCV*, 2019. 1, 2, 5, 6, 7, 12, 13, 17, 18
- [74] Songhyun Yu, Bumjun Park, Junwoo Park, and Jechang Jeong. Joint learning of blind video denoising and optical flow estimation. In *CVPRW*, 2020. 2
- [75] Feihu Zhang, Oliver J. Woodford, Victor Adrian Prisacariu, and Philip H.S. Torr. Separable flow: Learning motion cost volumes for optical flow estimation. In *ICCV*, 2021. 3, 11
- [76] Guozhen Zhang, Yuhan Zhu, Hongya Wang, Youxin Chen, Gangshan Wu, and Limin Wang. Extracting motion and appearance via inter-frame attention for efficient video frame interpolation. In *CVPR*, 2023. 2, 6
- [77] Haoxian Zhang, Yang Zhao, and Ronggang Wang. A flexible recurrent residual pyramid network for video frame interpolation. In *ECCV*, 2020. 1
- [78] Shiyu Zhao, Long Zhao, Zhixing Zhang, Enyu Zhou, and Dimitris Metaxas. Global matching with overlapping attention for optical flow estimation. In *CVPR*, 2022. 1
- [79] Tinghui Zhou, Shubham Tulsiani, Weilun Sun, Jitendra Malik, and Alexei A Efros. View synthesis by appearance flow. In *ECCV*, 2016. 1

Appendix

A. Architecture Details

We build three models with different sizes, termed AMT-S, AMT-L, and AMT-G. For reproducibility, the architecture details of them are shown in Fig. 9, Fig. 10, and Fig. 11, respectively. We employ standard residual blocks [12] and instance normalization [64] in the correlation encoder. The lookup radius is set to 3. For each update block, a bilinear upsampling layer follows each head on upper levels (*i.e.*, $l > 1$). The IFRBlock represents the decoder proposed in IFRNet [28], which jointly estimates the bilateral flows and the intermediate feature. To further improve performance, we upsample the correlation feature in the case of AMT-G to align its spatial resolution with the current interpolated feature, facilitating updates in the high-resolution space. The code is available at <https://github.com/MCG-NKU/AMT>.

B. Multi-Frame Interpolation

For the multi-frame setting, we use GoPro dataset [40] for training and evaluate our model on the test partition of GoPro dataset [40] and Adobe240 dataset [59]. Here, we aim at $8\times$ interpolation, synthesizing 7 intermediate frames with two input frames. The other training settings and loss functions are consistent with those in our main paper. Following recent frame interpolation works [18, 28], we inject a temporal embedding vector into the network for $8\times$ interpolation. The elements in this vector are all set to t according to the current time step, where $t \in \{1/8, 2/8, \dots, 7/8\}$. We compare our AMT-S with DVF [35], SuperSloMo [23], DAIN [1], and IFRNet-B [28]. The results of $8\times$ interpolation are shown in Tab. 3. Our method obtains the best PSNR and SSIM results on both evaluation datasets, indicating the

Method	GoPro [40]		Adobe240 [59]	
	PSNR	SSIM	PSNR	SSIM
DVF [35]	21.94	0.776	28.23	0.896
SuperSloMo [23]	28.52	0.891	30.66	0.931
DAIN [1]	29.00	0.910	29.50	0.910
IFRNet-B [28]	29.97	0.922	31.93	0.936
AMT-S (Ours)	30.20	0.927	32.04	0.938

Table 3. Quantitative comparison for $8\times$ interpolation.

effectiveness of the proposed AMT for the task of multi-frame interpolation. Fig. 7 and Fig. 12 visually compare our method and IFRNet-B on the Adobe240 dataset. Here, we visualize the cases for $1/4$ and $1/2$ time steps. It can be seen that our method can generate more temporally consistent results with fewer artifacts and more clear edges.

C. Limitation

Although our method has shown remarkable performance, the 4D correlation volume computed from all pairs of pixels makes it hard to adapt to very high-resolution inputs under a resource-constrained environment. This is because the computational complexity of constructing correlation volumes is quadratic to the image resolution. The possible ways to alleviate this problem include computing each correlation value only when it is looked up [61] or factorizing the 4D correlation volume to two 3D correlation volumes [70].

D. Discussions with RAFT

Teed and Deng [61] proposed RAFT, which iteratively performs lookups on multi-scale 4D correlation volumes for updating flow fields. Given its impressive results, current state-of-the-art flow estimation methods [17, 24, 70, 71, 75] all derive from such architecture design. Besides, it inspires the development of stereo matching [33] and scene flow [62]. However, the RAFT-like design paradigm is not well investigated in frame interpolation.

To better model large motions for frame interpolation, we build AMT based on RAFT. However, AMT involves many novel and task-specific designs beyond it. To better illustrate our model, we detail the differences between our AMT and RAFT from the following perspectives:

Volume Design: RAFT constructs a unidirectional correlation volume because it only needs to predict the optical flow along one direction. For frame interpolation, we hope to model the dense correspondences on both directions for updating bilateral flows. We thus construct bidirectional correlation volumes. We have verified the effectiveness of the bidirectional correlation volumes in Tab. 2a of the main paper.

Context Encoder: In RAFT, the context encoder extracts the content feature only from the first input frame. Because of the characteristics of frame interpolation, in our AMT, the context encoder takes the image pair as input. It outputs the initial intermediate feature, the initial bilateral flows, and the pyramid features from the input pair. This design is also inspired by recent one-stage frame interpolation methods [28, 50].

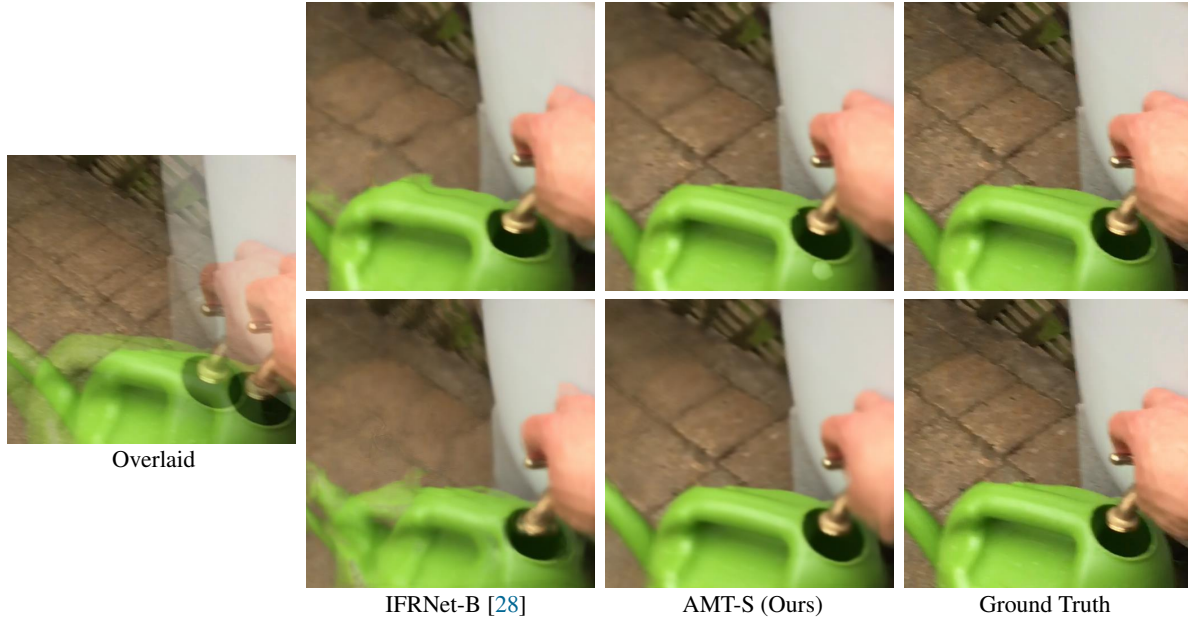


Figure 7. Qualitative results of our AMT-S and IFRNet-B [28] on Adobe240 [59]. The time steps are 1/4 and 1/2 from top to bottom.

Case	Vimeo90K [73]	SNU-FILM [8]		FLOPs (G)
		Hard	Extreme	
Single-scale Pred.	35.94	30.52	25.26	124
ConvGRU	35.99	30.58	25.27	132
Tied Weights	35.93	30.56	25.22	121
Convex Upsampling	35.99	30.56	25.28	123
Original Model	35.97	30.60	25.30	121

Table 4. Investigation on RAFT-like [61] designs. The default setting is marked in gray .

Correlation Lookup: The lookup operation can be directly performed in RAFT for the identical coordinate system between the correlation volume and predicted flow field. To solve the coordinate mismatch issue caused by the invisible frame, we propose to scale the bilateral flows before the lookup operation. Besides, we retrieve bidirectional correlations instead of the unidirectional ones in RAFT. We use the initial bilateral flows ($F_{t \rightarrow 0}^1, F_{t \rightarrow 1}^1$) as the initial starting point, while RAFT uses zero instead. The lookup strategy is investigated in Tab.2b of the main paper.

Predict and Update Manner: While RAFT predicts and updates the flow prediction at a single resolution, we predict and update the bilateral flows in a coarse-to-fine manner. We also provide a variant of our AMT to verify the design, which only predicts the flow fields at a single resolution before feeding into the last decoder. Tab. 4 shows that this variant performs worse than the original one. This indicates that predicting multi-scale flows are important for frame interpolation. Besides, we also investigate the effectiveness of the cross-scale update in Tab. 2d of the main paper.

Update Block: In the design of the update block, our AMT differs from RAFT in five aspects: 1) While RAFT regards the feature extracted from the visible frame as the content guidance, we use the interpolated intermediate feature representing the invisible frame

instead. 2) RAFT only has one head in update block for regressing a flow residual, while we have two heads for jointly predicting content and flow residuals. The two aspects mentioned above have been discussed in Tab. 2c of the main paper. 3) We stack two convolutional layers instead of a cumbersome ConvGRU unit in RAFT to handle the content and motion features. We also investigate a variant that equips with a ConvGRU unit in each update block. As shown in Tab. 4, this variant shows a comparable performance in contrast to the original one, but it has more computational costs. We thus choose to stack two convolutional layers for efficiency. 4) The weights of update blocks are not shared across levels in our AMT. However, weight tying is beneficial to RAFT. Tab. 4 demonstrates that the model with untied weights performs better than that with tied weights. 5) We employ bilinear upsampling instead of convex sampling in RAFT for upscaling the flow fields. As shown in Tab. 4, the two upsampling operators have similar performance, but convex upsampling will incur more computation costs. Thus, we choose the bilinear upsampling in our AMT.

Final Objective: RAFT is designed for flow estimation and is optimized only with flow regression loss. However, our AMT is introduced for frame interpolation and is supervised with both task-oriented flow distillation loss and distortion-oriented content losses. We need to consider not only the fidelity of estimated flows but also the diversity for meeting the requirement of task-oriented flows. We thus output multiple flow pairs rather than a single flow field in RAFT. Besides, occlusion reasoning and residual hallucination also need to be considered for faithful content generation.

E. Discussion about Multi-Field Refinement

Some works [15, 46, 47] also attempt to predict multiple flow pairs for preparing intermediate content candidates. Specifically, BMBC [46] predicts six bilateral motions through the bilateral motion network and optical flow approximation. ABME [47] gen-

erates four bilateral flow fields based on asymmetric motion assumption. After obtaining warped candidate frames and context features, the two works rely on a dynamic filter and even a cumbersome synthesis network to generate the final intermediate frame. Thus, they are inefficient for practical usage. In contrast, our AMT is more efficient, as shown in Tab. 1 of the main paper. We generate multiple flow fields in a single forward pass instead of multiple inference steps in BMBC and ABME. Besides, we obtain the intermediate candidates only in the image domain rather than the feature domain and stack two lightweight convolutional layers for fusing these candidates.

M2M-VFI [15] is most relevant to our multi-field refinement. It also generates multiple flows in one step and prepares warped candidates in the image domain. However, there are five key differences between our multi-field refinement and M2M-VFI. First, our method generates the candidate frames by backward warping rather than forward warping in M2M-VFI. Second, while M2M-VFI predicts multiple flows to overcome the hole issue and artifacts in overlapped regions caused by forward warping, we aim to alleviate the ambiguity issue in the occluded areas and motion boundaries by enhancing the diversity of flows. Third, M2M-VFI needs to estimate bidirectional flows first through an off-the-shelf optical flow estimator and then predict multiple bilateral flows through a motion refinement network. On the contrary, we directly estimate multiple bilateral flows in a one-stage network. In this network, we first estimate one pair of bilateral flows at the coarse scale and then derive multiple groups of fine-grained bilateral flows from the coarse flow pairs. Fourth, M2M-VFI jointly estimates two reliability maps together with all pairs of bilateral flows, which can be further used to fuse the overlapping pixels caused by forward warping. As shown in Eqn. (5) of the main paper, we estimate not only an occlusion mask but a residual content for cooperating with each pair of bilateral flows. The residual content is used to compensate for the unreliable details after warping. This design has been investigated in Tab. 2e of the main paper. Fifth, we stack two convolutional layers to adaptively merge candidate frames, while M2M-VFI normalizes the sum of all candidate frames through a pre-computed weighting map.

F. More Visual Results

In this section, we provide additional visual results on two benchmark datasets, including Vimeo90K [73] and SNU-FILM [8], to further show the superiority of the proposed AMT. The comparison methods include CAIN [8], AdaCoF [29], ABME [47], RIFE [18], IFRNet(-B/-L) [28], and VFIFormer [37]. For a fair comparison, we also divide these methods into two groups according to the computational cost. As shown in Fig. 8, 13-18, our AMT synthesizes the object with large motions more faithfully and generates plausible textures with fewer artifacts.

G. Broader Impact

As presented in this paper, our AMT can synthesize faithful non-existent frames between two visible frames. Given its reliable synthesis results, our method may be abused to forge or tamper with videos.



Figure 8. Qualitative comparison between AMT-G with VFIFormer. Our method recovers more clear structure and edges.

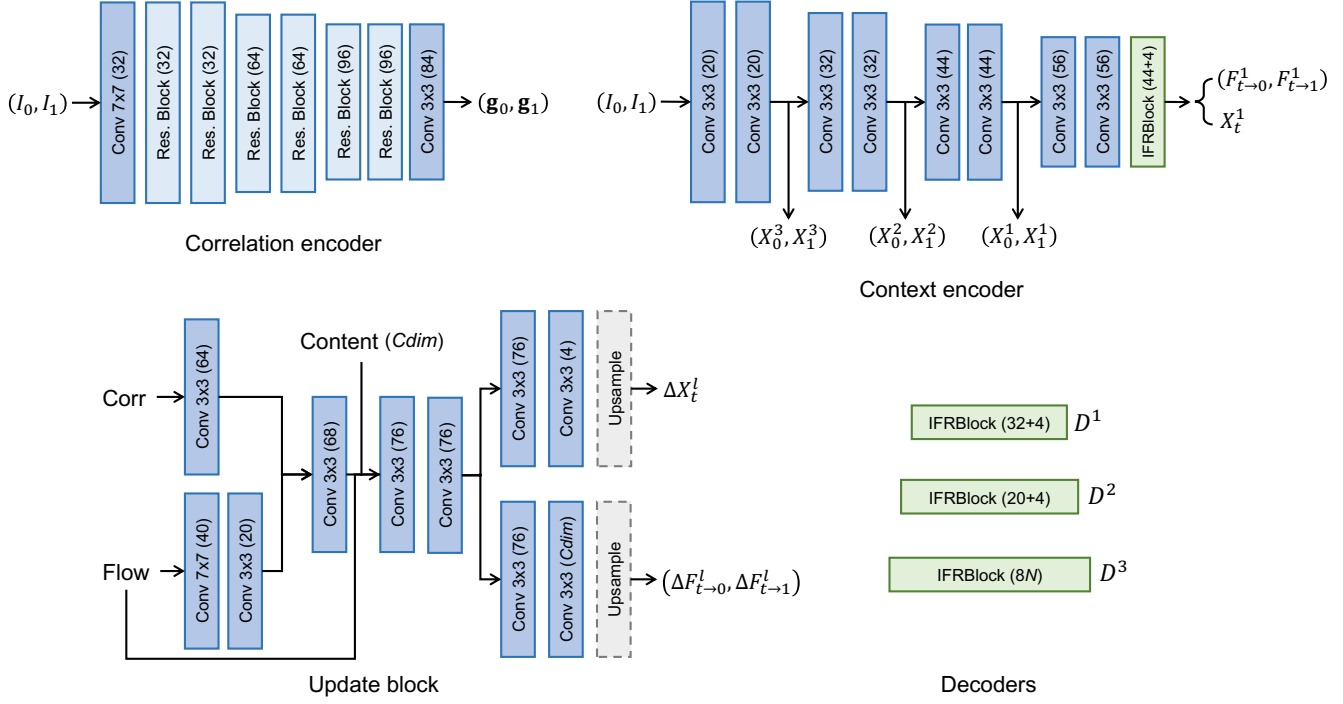


Figure 9. Architecture details of the AMT-S. The number in parentheses denotes the output channels. N represents the number of output groups. IFRBlock denotes the decoder proposed in IFRNet [28].

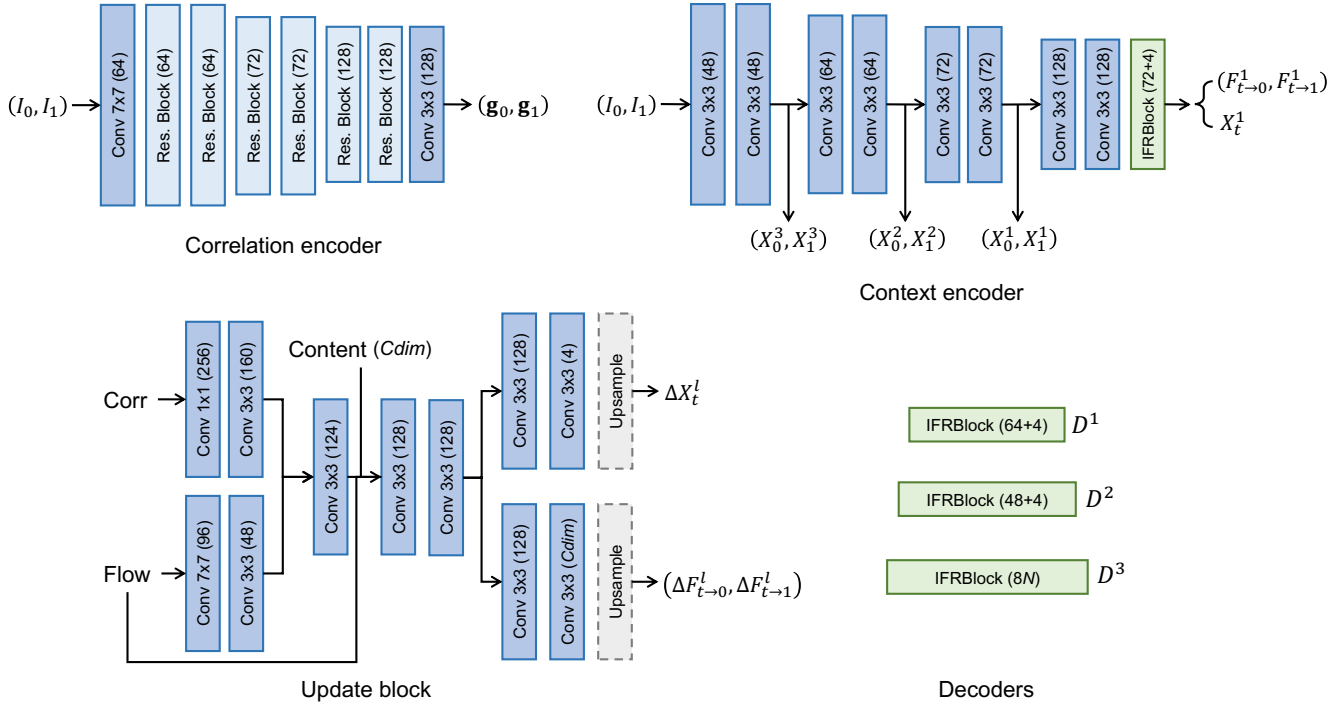


Figure 10. Architecture details of the AMT-L. The number in parentheses denotes the output channels. N represents the number of output groups. IFRBlock denotes the decoder proposed in IFRNet [28].

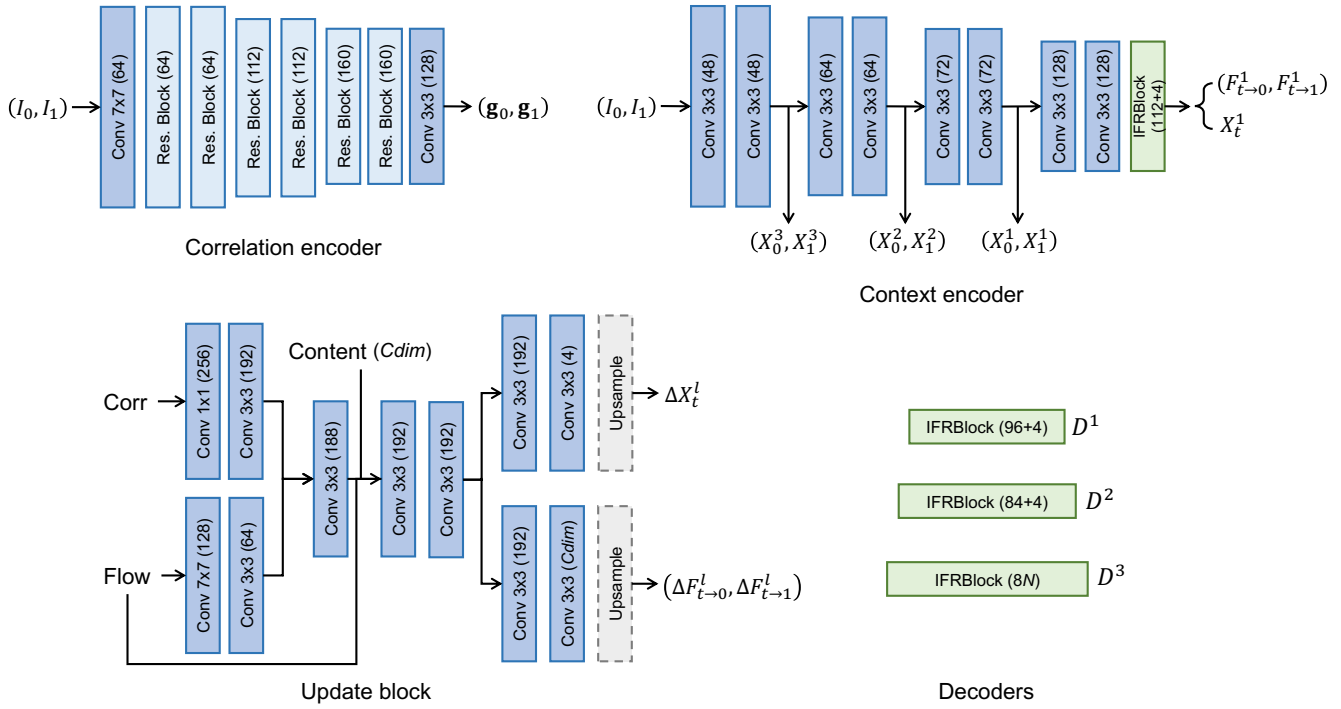


Figure 11. Architecture details of the AMT-G. The number in parentheses denotes the output channels. N represents the number of output groups. IFRBlock denotes the decoder proposed in IFRNet [28].

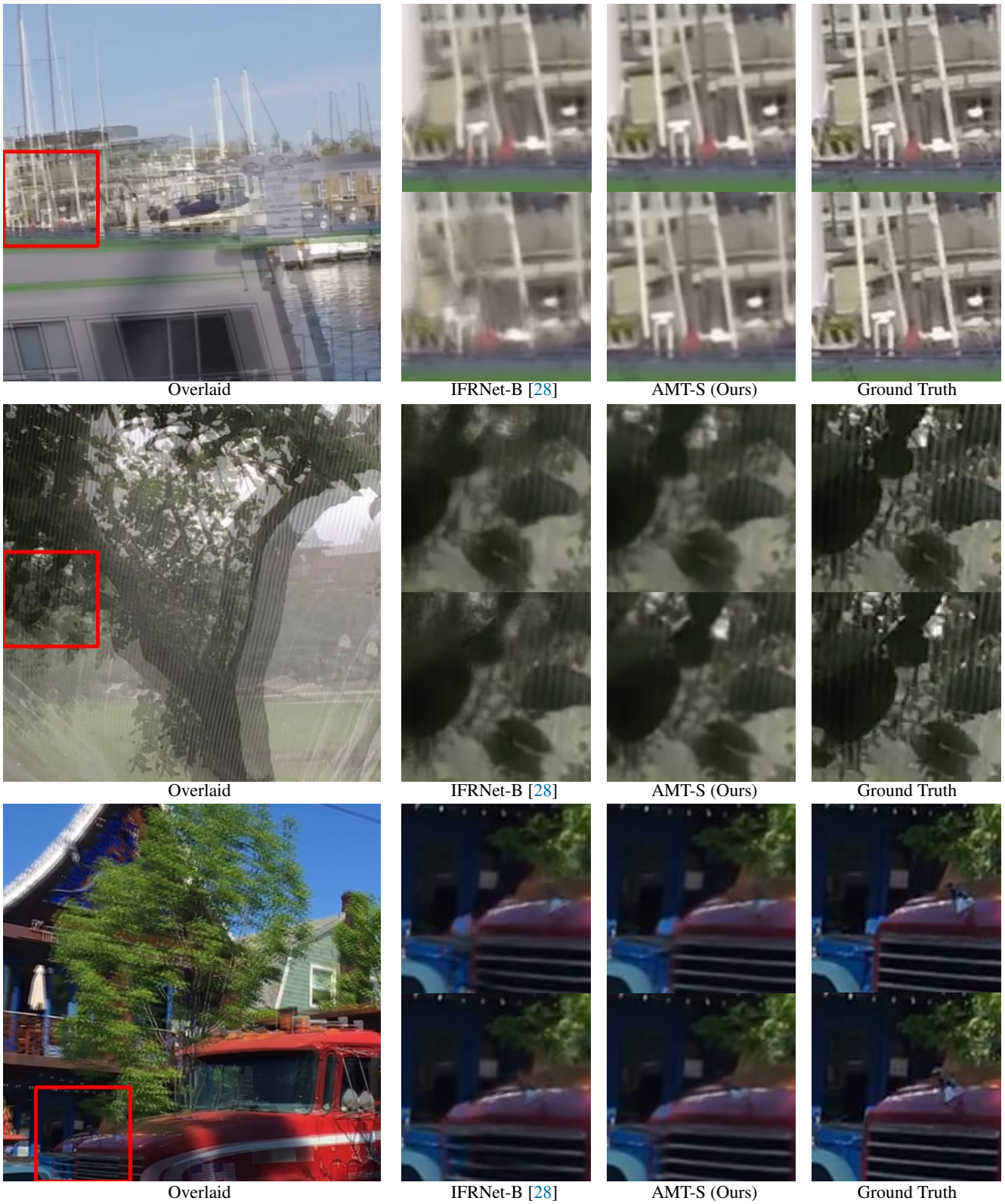


Figure 12. Qualitative results of AMT-S and IFRNet-B [28] on Adobe240 [59]. The time steps are 1/4 and 1/2 from top to bottom.

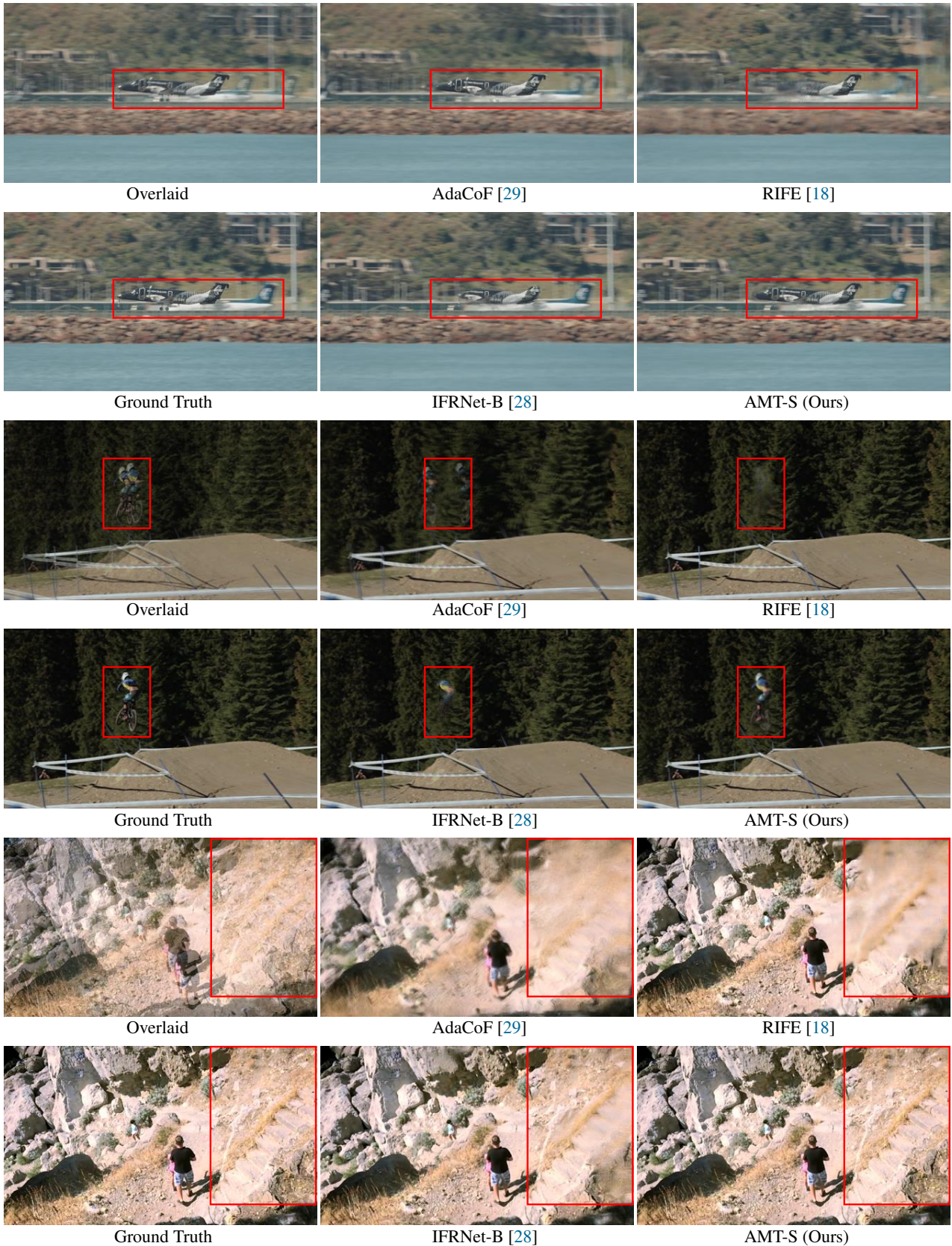


Figure 13. Visual comparison for the methods with low computational complexity on Vimeo90K dataset [73]

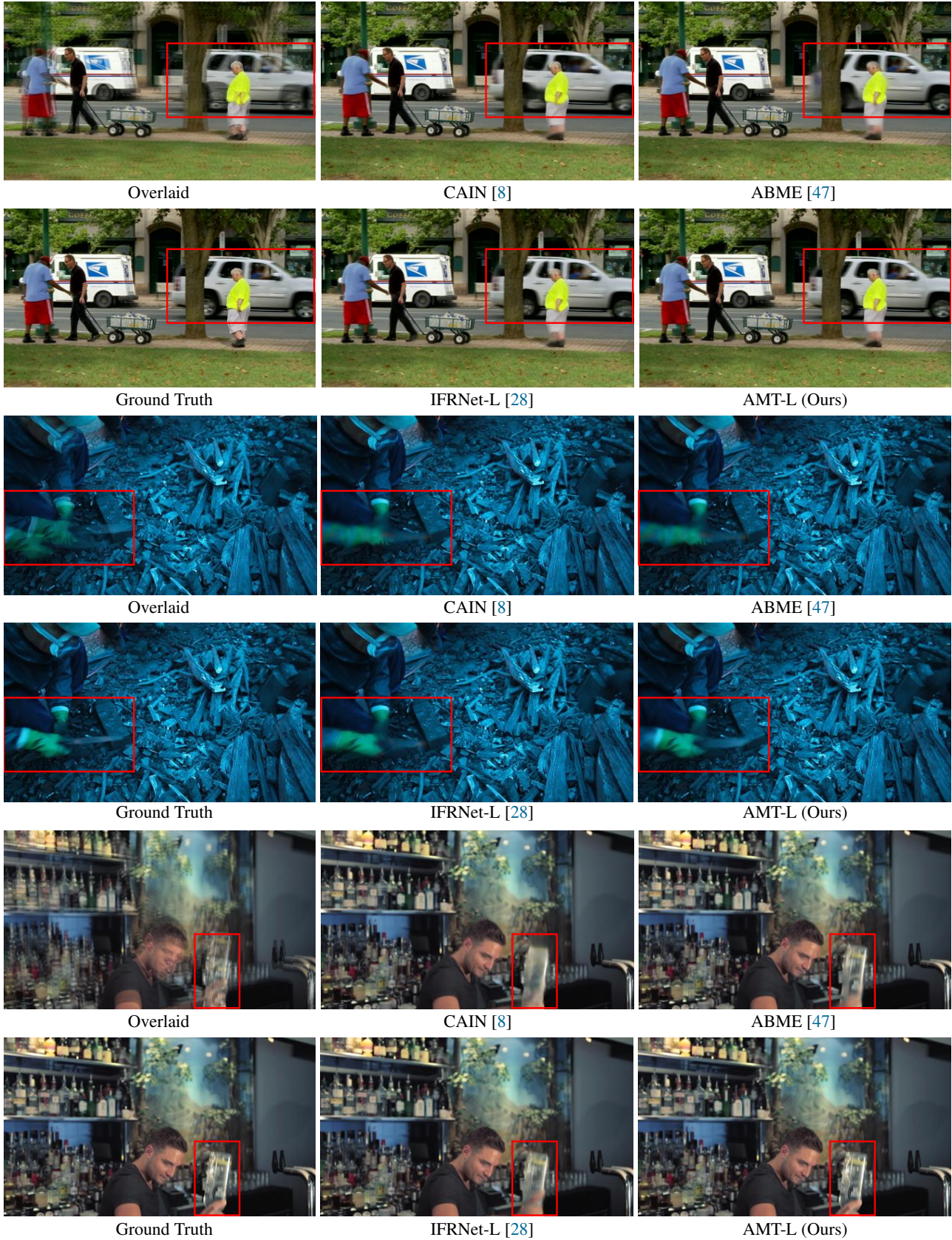


Figure 14. Visual comparison for the methods with relatively high computational complexity on Vimeo90K dataset [73].

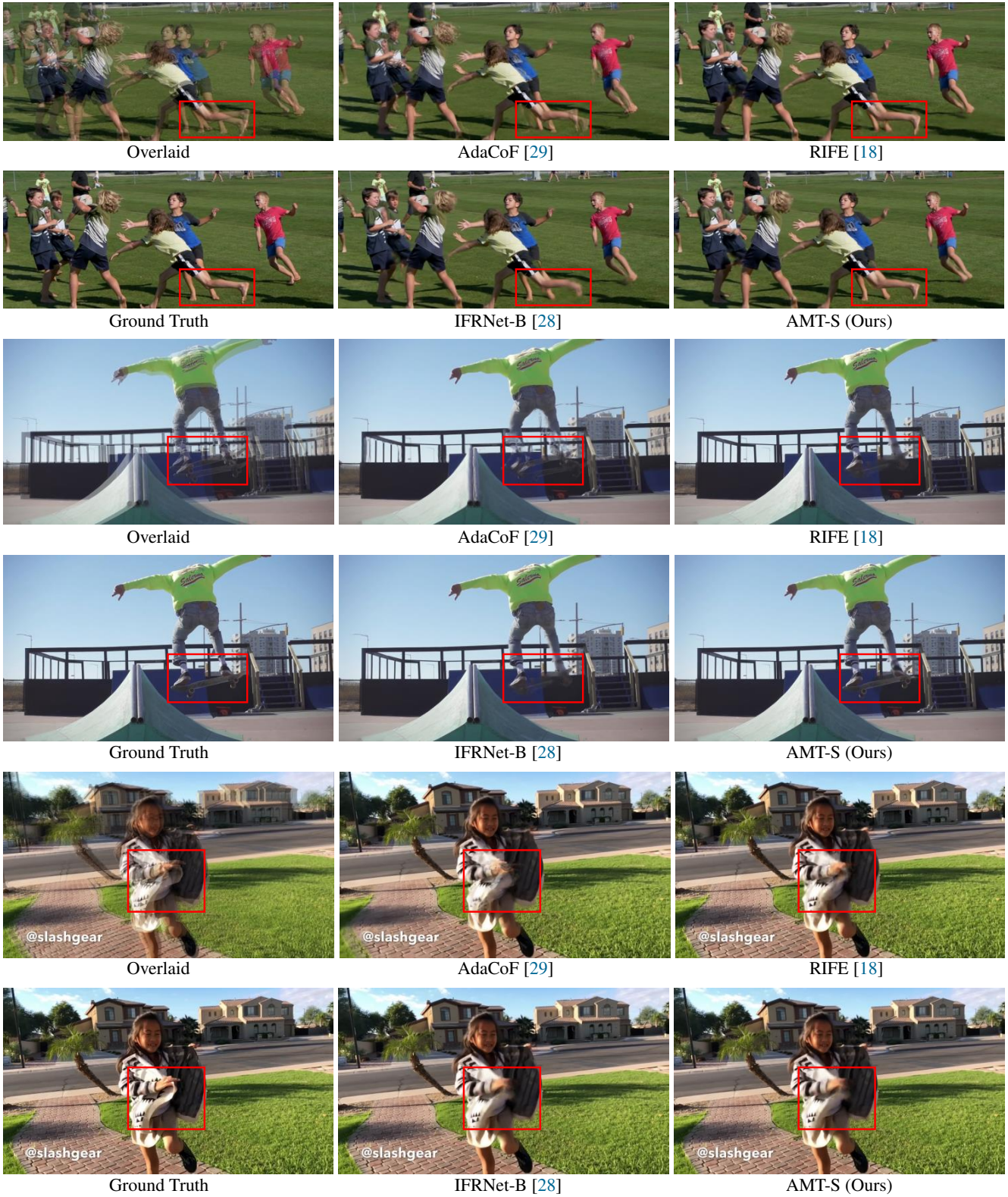


Figure 15. Visual comparison for the methods with low computational complexity on the Hard partition in SNU-FILM dataset [8].

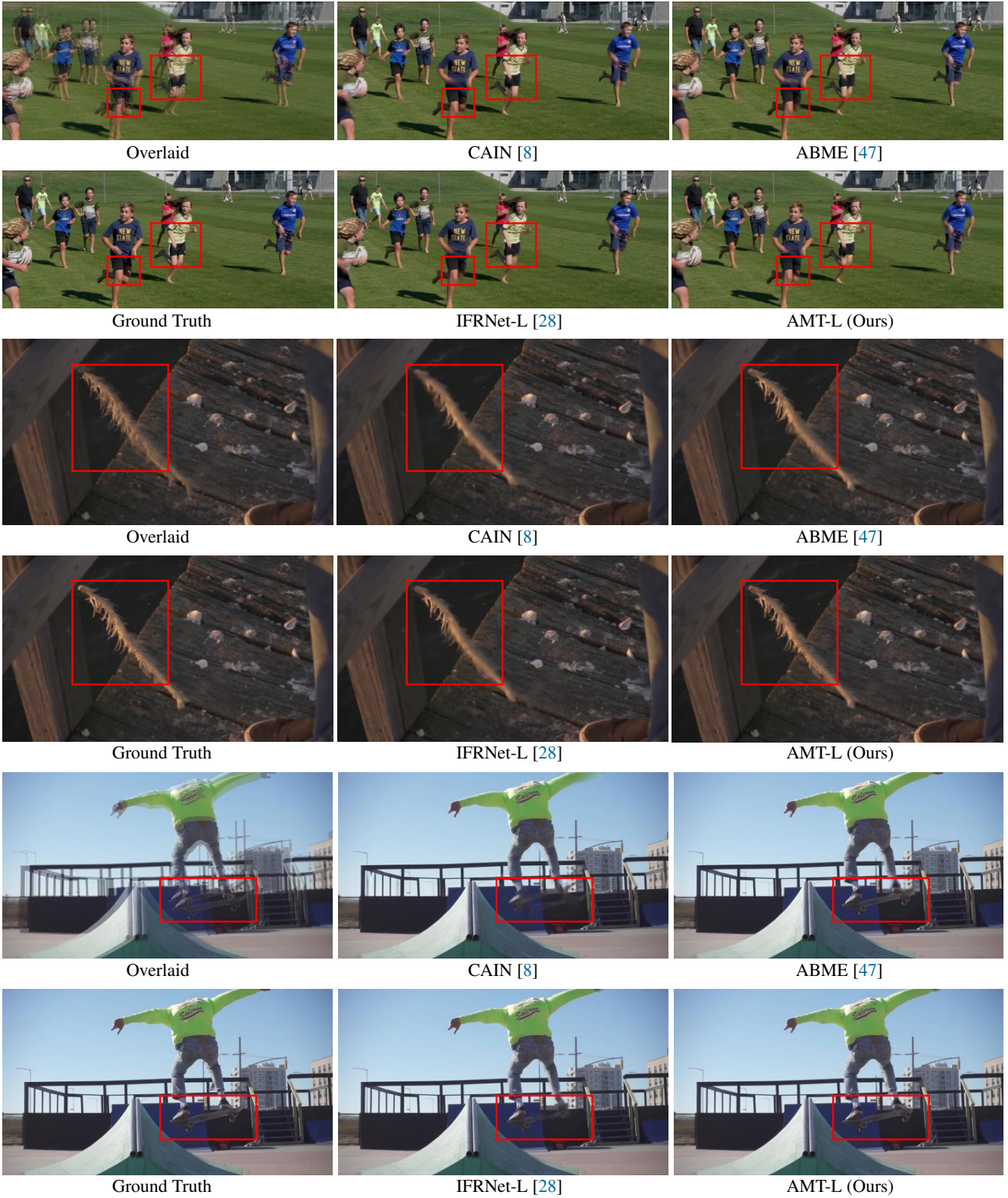


Figure 16. Visual comparison for the methods with relatively high computational complexity on the Hard partition in SNU-FILM dataset [8].

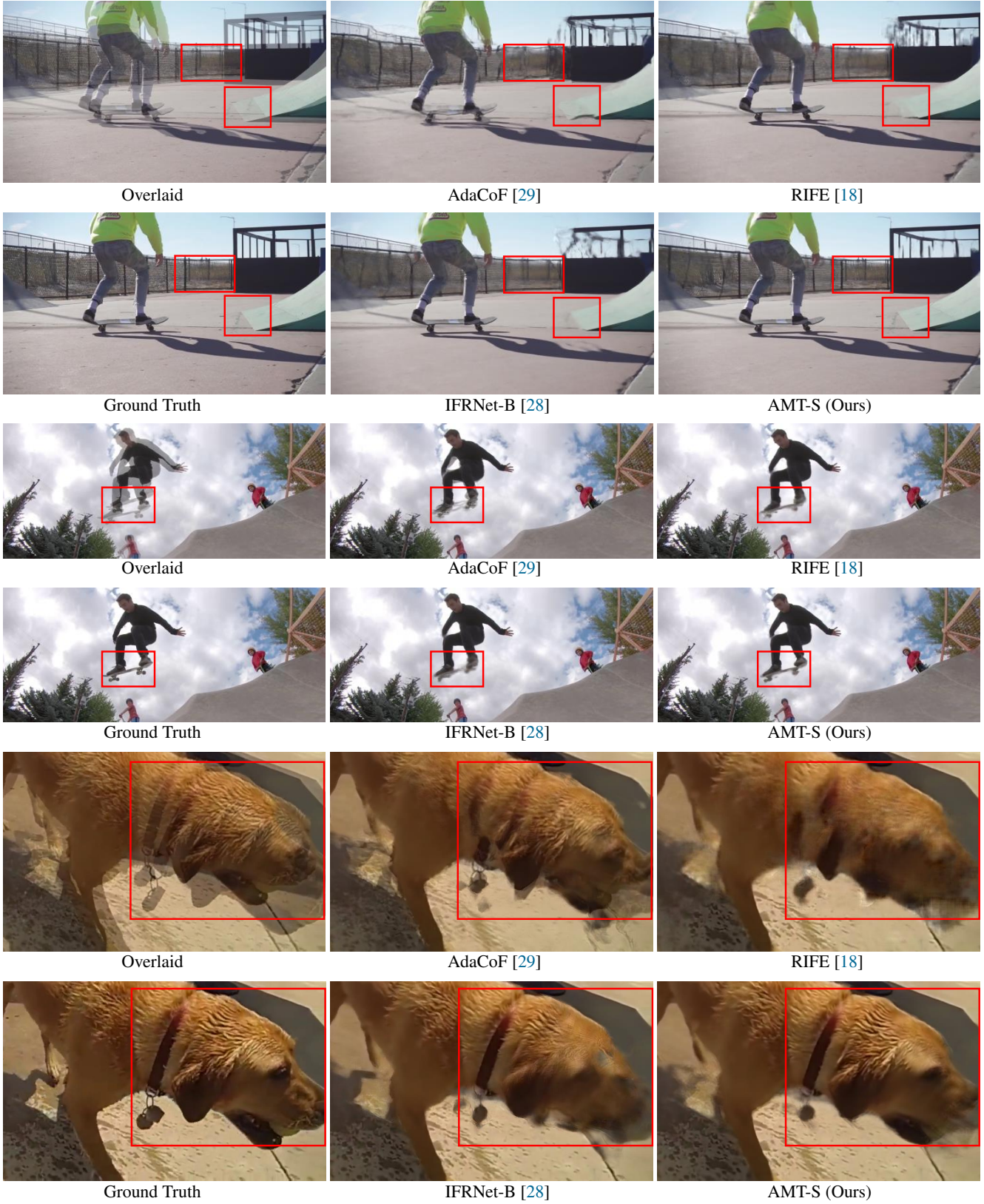


Figure 17. Visual comparison for the methods with low computational complexity on the Extreme partition in SNU-FILM dataset [8].

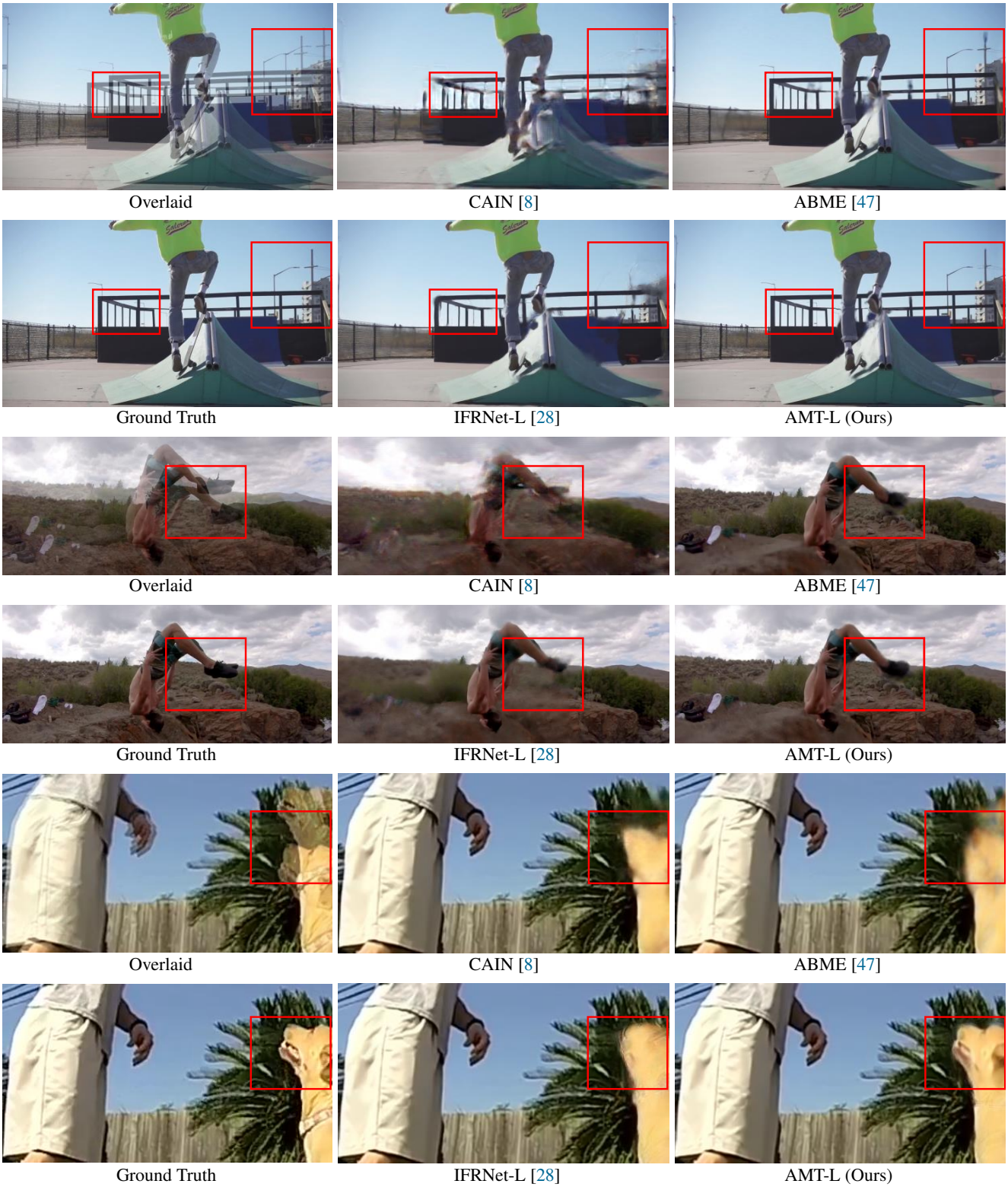


Figure 18. Visual comparison for the methods with relatively high computational complexity on the Extreme partition in SNU-FILM dataset [8].