

# Paper Reading

Dachun Kai

USTC

October 18, 2021

## XVFI: eXtreme Video Frame Interpolation

Hyeonjun Sim\*

Jihyong Oh\*

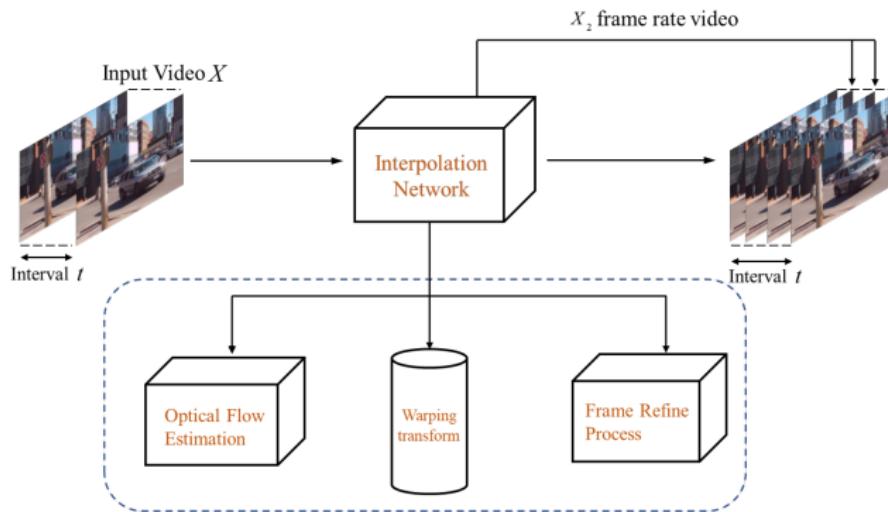
Munchurl Kim<sup>†</sup>

Korea Advanced Institute of Science and Technology

{f1hy5836, jhoh94, mkimee}@kaist.ac.kr

# Task: Video Frame Interpolation

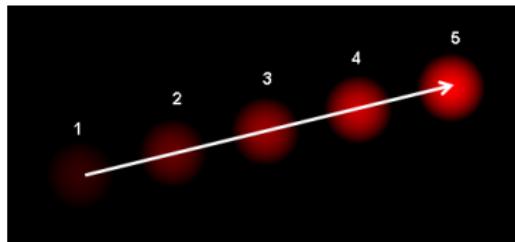
- A **low-level** vision task to increase frame rate of a video.
- Method: **motion-based**, kernel-based, phase-based, **event-based**
- **Flow-based** method framework:



- Evaluation metric: *PSNR*, *SSIM*, *IE*(Interpolation Error)

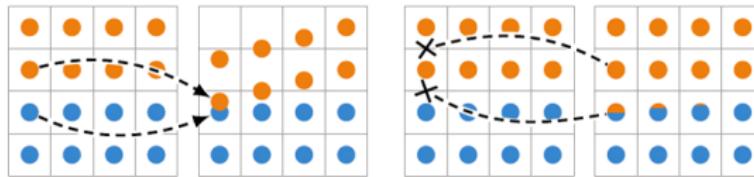
# Preliminaries: Optical Flow and Backward Warping

- Optical Flow depicts **motion** of objects in a visual scene.

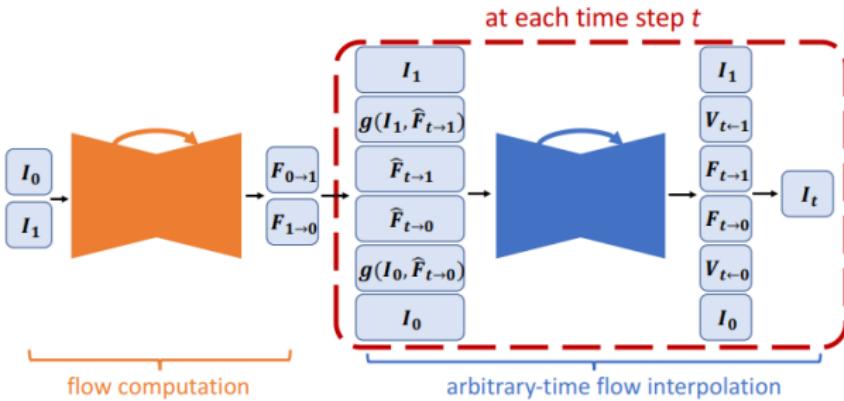


- Flow Example:

- Image Warping



# Milestone: Super SloMo<sup>1</sup>



## Ideas

$$\hat{F}_{t \rightarrow 0} = -(1-t)tF_{0 \rightarrow 1} + t^2F_{1 \rightarrow 0}$$

$$\hat{F}_{t \rightarrow 1} = (1-t)^2F_{0 \rightarrow 1} - t(1-t)F_{1 \rightarrow 0}$$

$$\hat{I}_t = \frac{1}{Z} \odot ((1-t)V_{t \leftarrow 0} \odot g(I_0, F_{t \rightarrow 0}) + tV_{t \leftarrow 1} \odot g(I_1, F_{t \rightarrow 1}))$$

<sup>1</sup>Huaizu Jiang et al. "Super slomo: High quality estimation of multiple intermediate frames for video interpolation". In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. 2018, pp. 9000–9008.

# Problem Definition

- Motivation

- ▶ Handle the VFI for 4K videos with large motion, computed by IRR<sup>2</sup>.

---

<sup>2</sup>Junhwa Hur and Stefan Roth. "Iterative residual refinement for joint optical flow and occlusion estimation". In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* 2019, pp. 5754–5763.

# Proposed X4K1000FPS Dataset

- Impressive *4K@1000fps* dataset.



- Various objects: crowds, cars, trains, plants, boats...
- Various places: stadiums, stations, beaches, rivers...

# Proposed X4K1000FPS Dataset

- Impressive *4K@1000fps* dataset.

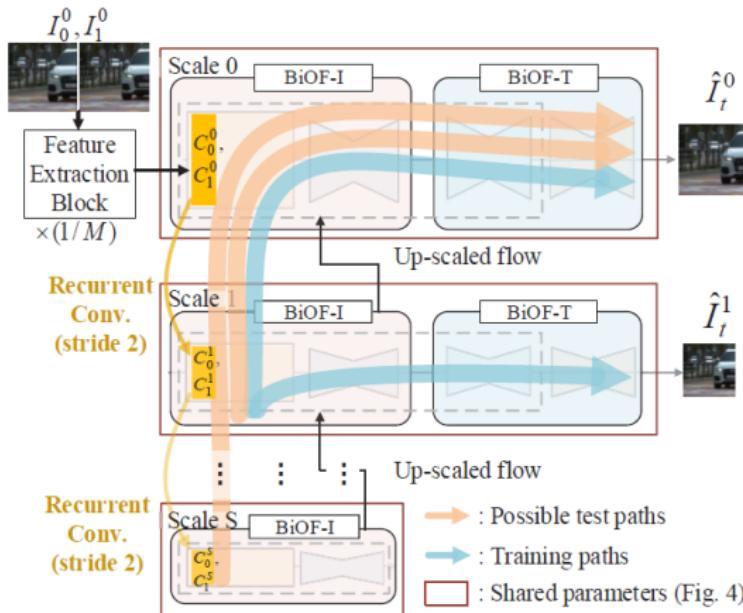
Dataset	Occlusion [16]			Flow magnitude [16]		
	25 <sub>th</sub>	50 <sub>th</sub>	75 <sub>th</sub>	25 <sub>th</sub>	50 <sub>th</sub>	75 <sub>th</sub>
Vimeo90K [48]	6.8	11.9	18.1	3.1	4.9	7.1
Adobe240fps [39]	0.8	1.7	3.2	3.8	8.9	16.3
X-TEST (ours)	2.1	5.6	17.7	23.9	81.9	138.5
X-TRAIN (ours)	6.9	10.1	15.7	5.5	18.0	59.5

25<sub>th</sub>, 50<sub>th</sub> and 75<sub>th</sub> represent percentiles of each dataset.

- Detailedly analysis: occlusion, flow magnitude
- Select extreme scenes for X-TRAIN and X-TEST.

# Proposed Method

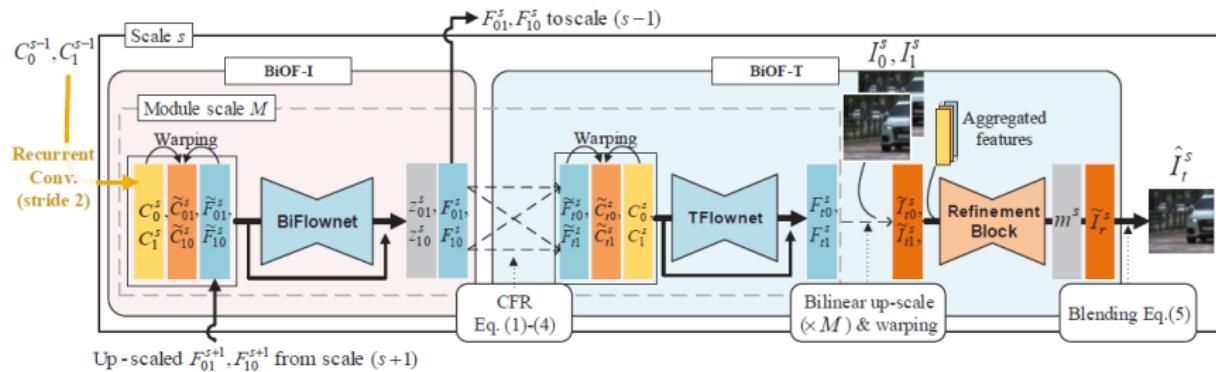
- Feature Extraction: strided conv + residual block



- Scale Adaptivity: **Shared** parameters, adapt to resolution of  $I_0, I_1$

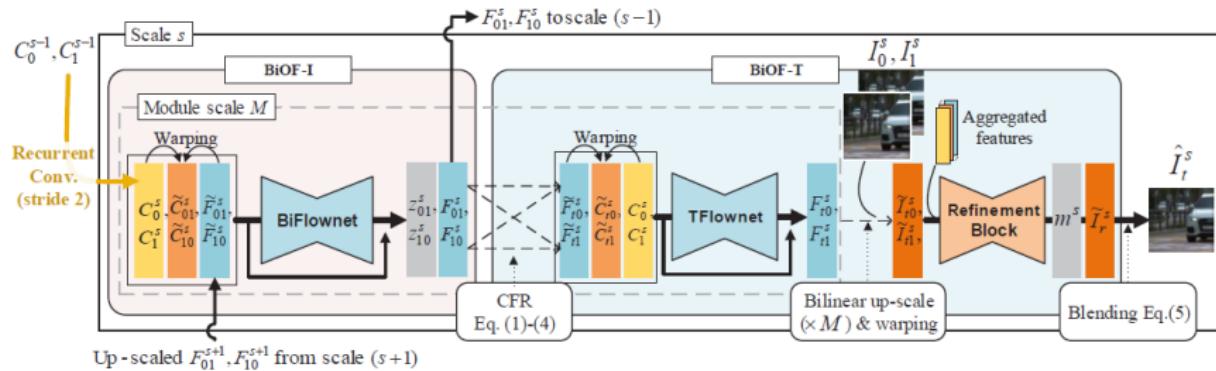
# Proposed Method

- Arch of XVFI-Net in scale  $s$ .



# Proposed Method

- Arch of XVFI-Net in scale  $s$ .

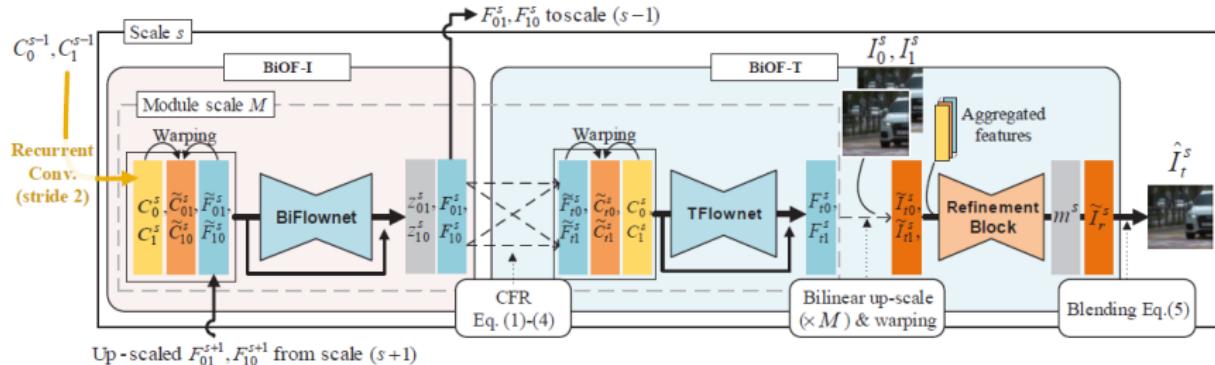


- BIOF-I module

- ▶ Input: context info  $C_0^{s-1}, C_1^{s-1}$ ,  $s + 1$ -scale Flow  $F_{01}^{s+1}, F_{10}^{s+1}$
- ▶ Output:  $s$ -scale Flow  $F_{01}^s, F_{10}^s$

# Proposed Method

- Arch of XVFI-Net in scale  $s$ .



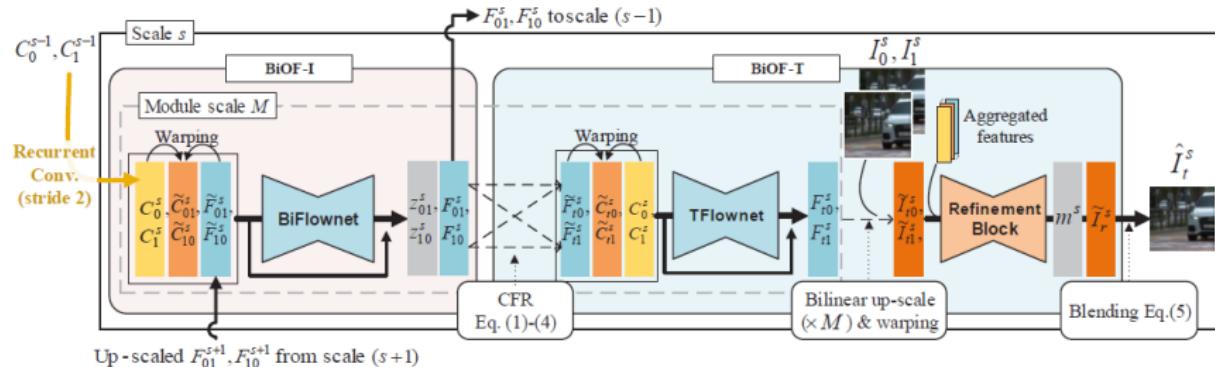
- BIOF-T module:

$$\tilde{F}_{t0}^x = \frac{(1-t) \sum_{\mathcal{N}_0} w_0 \cdot (-F_{0t}^y) + t \sum_{\mathcal{N}_1} w_1 \cdot F_{1 \cdot (1-t)}^y}{(1-t) \sum_{\mathcal{N}_0} w_0 + t \sum_{\mathcal{N}_1} w_1} \quad (1)$$

$$\tilde{F}_{t1}^x = \frac{(1-t) \sum_{\mathcal{N}_0} w_0 \cdot F_{0 \cdot (1-t)}^y + t \sum_{\mathcal{N}_1} w_1 \cdot (-F_{1t}^y)}{(1-t) \sum_{\mathcal{N}_0} w_0 + t \sum_{\mathcal{N}_1} w_1} \quad (2)$$

# Proposed Method

- Arch of XVFI-Net in scale  $s$ .



- BIOF-T module:

$$\mathcal{N}_0 = \{y \mid \text{round } (y + F_{0t}^y) = x\} \quad (3)$$

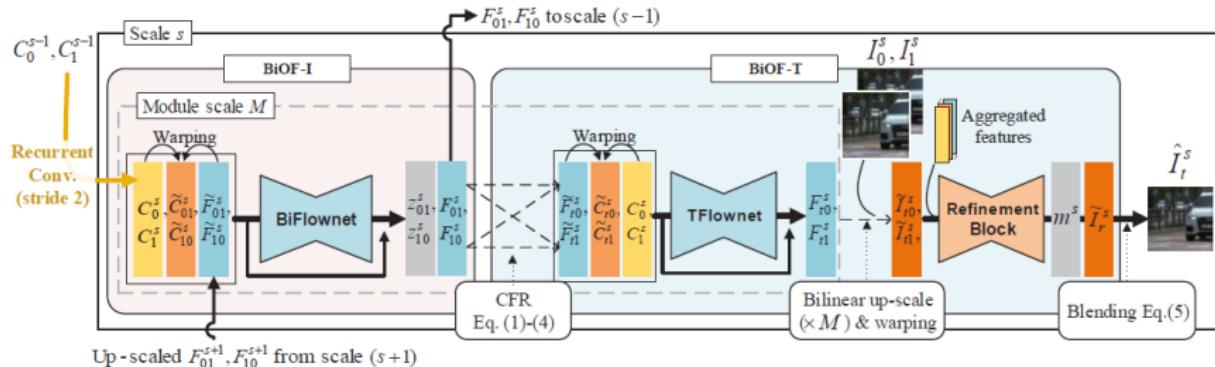
$$\mathcal{N}_1 = \{y \mid \text{round } (y + F_{1t}^y) = x\} \quad (4)$$

$$w_i = z_i^y \cdot G(|x - (y + F_{it}^y)|) \quad (5)$$

$$G(d) = e^{-d^2/\sigma^2} \quad (6)$$

# Proposed Method

- Arch of XVFI-Net in scale  $s$ .

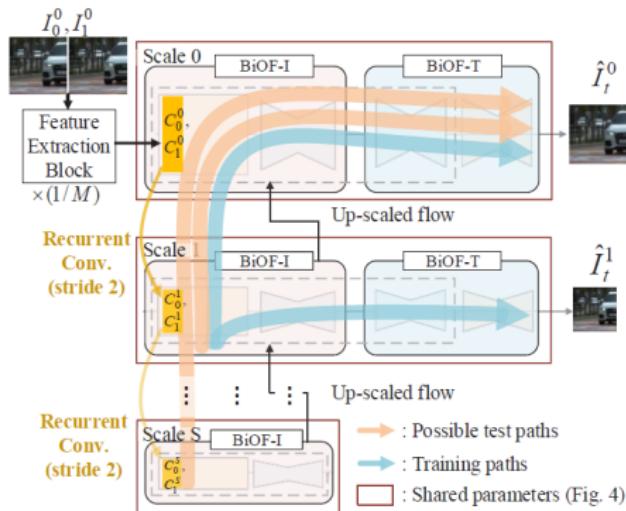


- BiOF-T module:

$$\hat{I}_t^s = \frac{(1-t) \cdot m^s \cdot \tilde{I}_{t0}^s + t \cdot (1-m^s) \cdot \tilde{I}_{t1}^s}{(1-t) \cdot m^s + t \cdot (1-m^s)} + \tilde{I}_r^s$$

where  $m^s$  is occlusion mask and  $\tilde{I}_r^s$  is residual image.

# Proposed Method

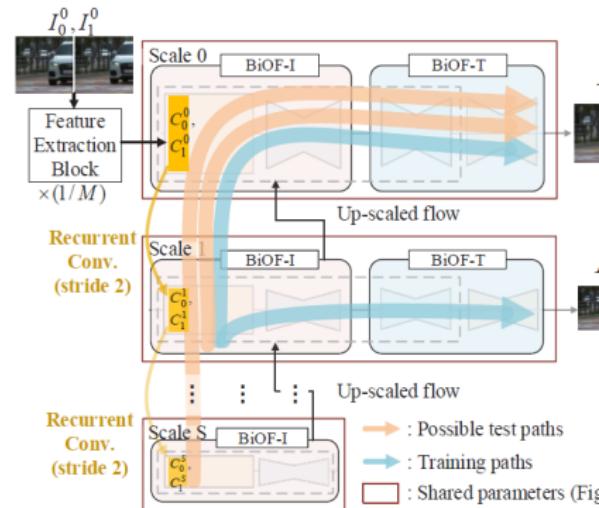


## Loss Function

- $\mathcal{L}_{\text{total}} = \mathcal{L}_r + \lambda_s \cdot \mathcal{L}_s$
- Multi-scale reconstruction loss:  $\mathcal{L}_r = \sum_{s=0}^{S_{tm}} \left\| \hat{I}_t^s - I_t^s \right\|_1$
- Edge-aware smoothness loss:  $\mathcal{L}_s = \sum_{i=0,1} \exp(-e^2 \sum_c |\nabla_x I_{tc}^0|)^\top \cdot |\nabla_x F_{ti}^0|$

# Model Analysis

- Shared parameters → **lightweighted**.
- Inference can start from any scale level.
- Coarse-to-fine strategy, but start from **low spatial resolution**.
- Interpolate a frame at **any time** instance from two input frames.



# Comparison to SOTA

- Settings

- ▶ Retrain SOTA on X-TRAIN with sub  $f$ , original with sub  $o$ .
- ▶ XVFI-Net train( $S_{trn} = 3$ ) and test( $S_{tst} = 3$  or 5).

- Results:

Methods ( $\times N$ )	X-TEST (PSNR/SSIM/tOF)	Adobe240fps (PSNR/SSIM)	#P ↓	R <sub>t</sub> ↓
AdaCoF <sub>o</sub> ( $\times 5.8$ )	23.90/0.727/6.89	25.26/0.785	<u>21.8</u>	<b>0.005</b>
AdaCoF <sub>f</sub> [25]	25.81/0.772/6.42	25.21/0.791	<u>21.8</u>	<b>0.005</b>
FeFlow <sub>o</sub> ( $\times 5.3$ )	24.00/0.756/6.59	25.18/0.785	102.5	1.681
FeFlow <sub>f</sub> [13]	25.16/0.783/6.54	24.17/0.780	102.5	1.681
DAIN <sub>o</sub> ( $\times 9.3$ )	26.78/0.807/3.83	29.89/ <u>0.911</u>	24	1.375
DAIN <sub>f</sub> [4]	27.52/0.821/3.47	29.99/0.910	24	1.375
Ours ( $S_{tst}=3$ )	<u>28.86/0.858/2.67</u>	<b>30.29/0.912</b>	<b>5.5</b>	<u>0.074</u>
Ours ( $S_{tst}=5$ )	<b>30.12/0.870/2.15</b>	<u>30.18/0.911</u>	<b>5.5</b>	0.075

$\times N$ : The ratio of number of iterations of the original version to that of - retrained version in the fair condition. #P: The number of parameters (M). R<sub>t</sub>: The runtime on 1024×1024-sized frames in sec.

**RED**: Best performance, BLUE: Second best performance.

Table 2. Quantitative comparisons on both X-TEST (4K) and Adobe240fps (HD) [39] for multi-frame interpolation ( $\times 8$ ).

# Comparison to SOTA

- Settings

- ▶ Retrain SOTA on X-TRAIN with sub  $f$ , original with sub  $o$ .
- ▶ XVFI-Net train( $S_{trn} = 3$ ) and test( $S_{tst} = 3$  or  $5$ ).

- Results:

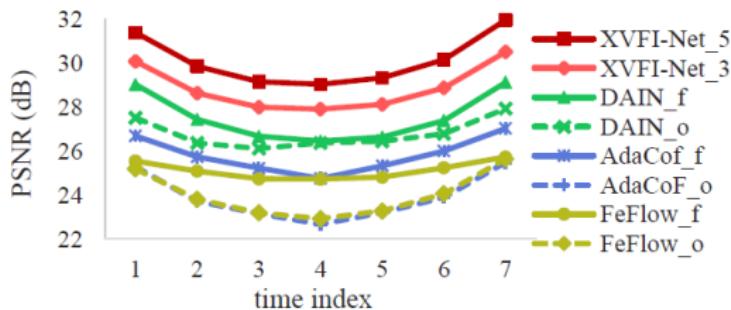


Figure 5. PSNR profiles for multi-frame interpolation results ( $\times 8$ ) on X-TEST.

# Comparison to SOTA

- Qualitative Results:

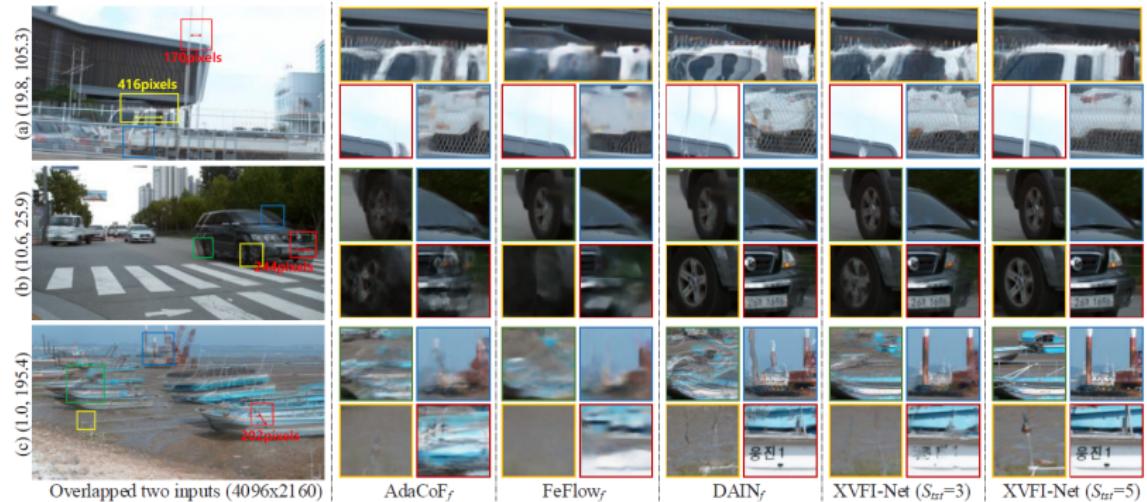
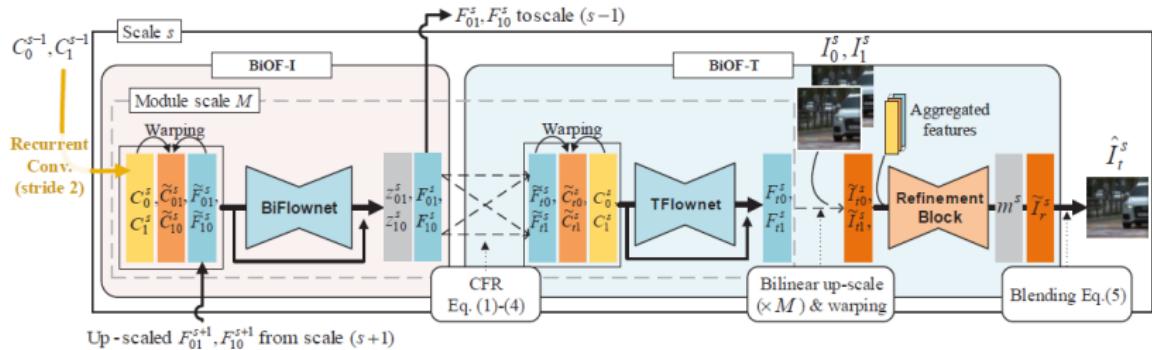


Figure 6. Visual comparisons for VFI results ( $t = 0.5$ ) on X-TEST for our and retrained SOTA methods with X-TRAIN. (\*, \*): occlusions and optical flow magnitudes between the two input frames measured by [16]. *Best viewed in zoom.*

# Ablation Studies



# Ablation Studies

- Flow Approximation

Methods \ Metrics	EPE↓	PSNR↑	SSIM↑	tOF↓
(a) Linear comb.	0.0752	28.73	0.8518	2.89
(b) Flow reversal	0.0892	28.30	0.8425	2.98
(c) CFR (ours)	<b>0.0721</b>	<b>28.86</b>	<b>0.8582</b>	<b>2.67</b>

- Adjustable scalability

$S_{trn}$ \ $S_{tst}$	(PSNR(dB)↑ / SSIM[44]↑ / tOF[8]↓)		
	1	3	5
1	26.85/0.806/4.90	<b>28.40/0.852/3.46</b>	27.14/0.842/3.69
3	23.61/0.729/6.56	29.22/0.863/2.68	<b>30.35/0.879/1.98</b>
5	22.37/0.699/6.71	23.70/0.724/6.39	<b>29.48/0.864/2.08</b>

**RED:** Best performance of each row

# Ablation Studies

- Robustness on LR-LFR dataset

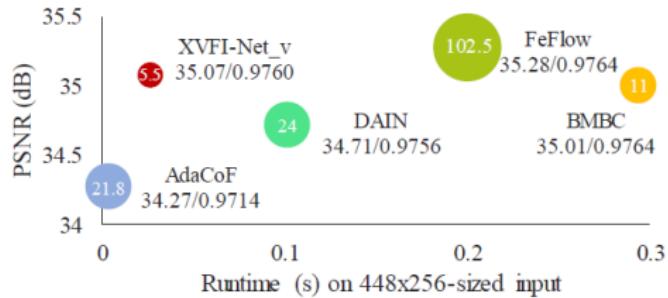


Figure 8. PSNR/SSIM vs runtime (s) on Vimeo90K [48] with model size (M) indicated in each circle.

# Failure cases

- Cases on 4K patches, with very large flow magnitude(196.5)

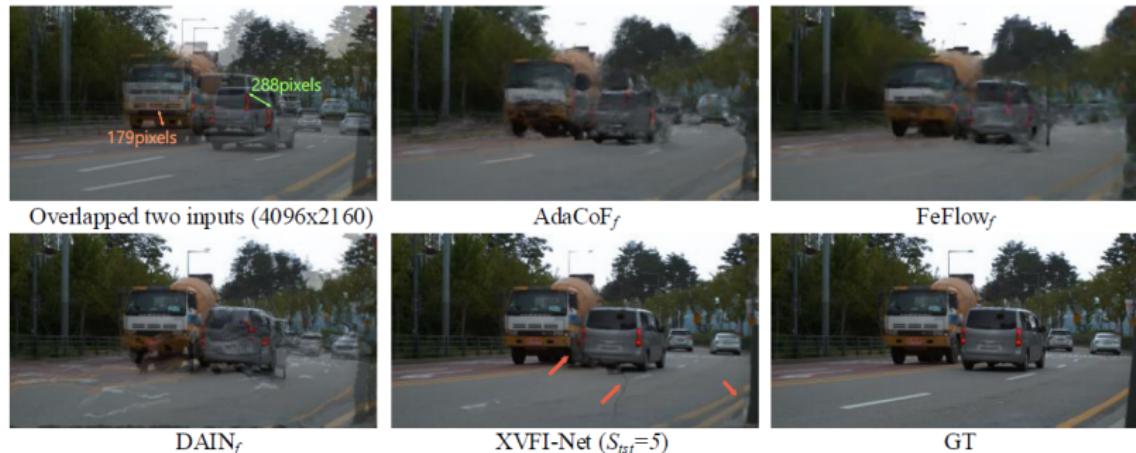


Figure 6. Failure cases of 4K result ( $t = 0.5$ ) on X-TEST for our and *retrained SOTA* methods with X-TRAIN, including the corresponding ground truth. *Best viewed in zoom.*

# Failure cases

- Cases on cropped patches



Figure 7. Failure cases of cropped results ( $t = 0.5$ ) on X-TEST for our and *retrained* SOTA methods with X-TRAIN, including the corresponding ground truth. *Best viewed in zoom.*

- Analysis: attribute to flow and warping algorithm.

# Conclusion

- Contributed HR-HFR **X4K1000FPS** dataset is very **valuable** for research.
- XVFI-Net showed **SOTA** performance on HR dataset, also **robust** to LR-LFR dataset.
- Extend VFI for more recent **real-world applications** with HR video.

*Thanks, Q & A*