

# Paper Reading

Dachun Kai

USTC

December 17, 2021

## Training Weakly Supervised Video Frame Interpolation with Events

Zhiyang Yu<sup>† 1, 2</sup>, Yu Zhang<sup>\* 2, 3</sup>, Deyuan Liu<sup>† 2, 5</sup>, Dongqing Zou<sup>2, 4</sup>,  
Xijun Chen<sup>\* 1</sup>, Yebin Liu<sup>3</sup>, and Jimmy Ren<sup>2, 4</sup>

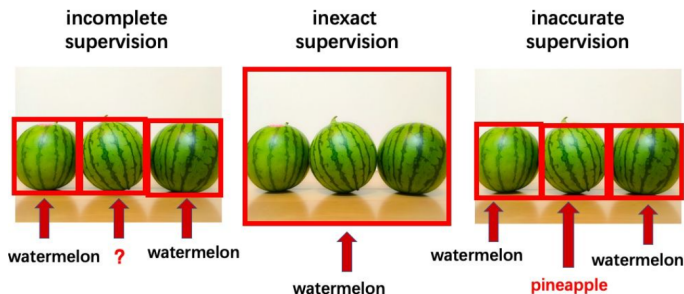
<sup>1</sup>Harbin Institute of Technology, <sup>2</sup>SenseTime Research and Tetras.AI, <sup>3</sup>Tsinghua University

<sup>4</sup>Qing Yuan Research Institute, Shanghai Jiao Tong University, <sup>5</sup>Peking University

# Preliminaries

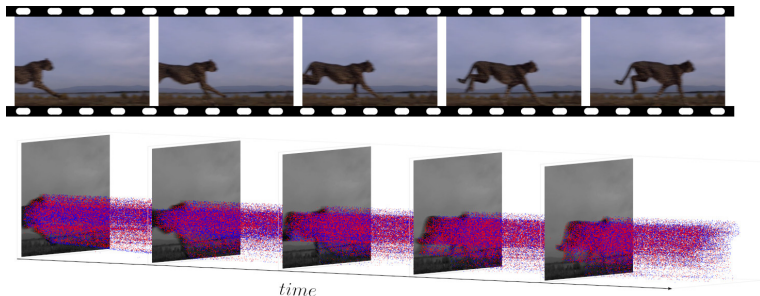
## Weakly Supervised learning

- Incomplete supervision: subset labels  $\rightarrow$  trained on low frame-rate videos
- Inexact supervision: coarse-grained labels
- Inaccurate supervision: wrong labels



# Preliminaries

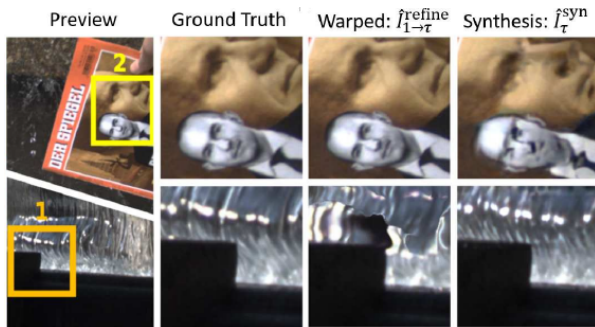
## Great performance of Event Camera



- High temporal resolution(microsecs)
- Low Latency
- Low Power(SNN, neuromorphic chip)
- No motion blur
- High dynamic range

# Motivation

- Parameterized **motion model** falls into an ill-posed problem when input frames are **sparse**.
- Event-based methods: learn **frame residual**, can't solve **low-texture surface**, always need high frame-rate auxiliary videos.

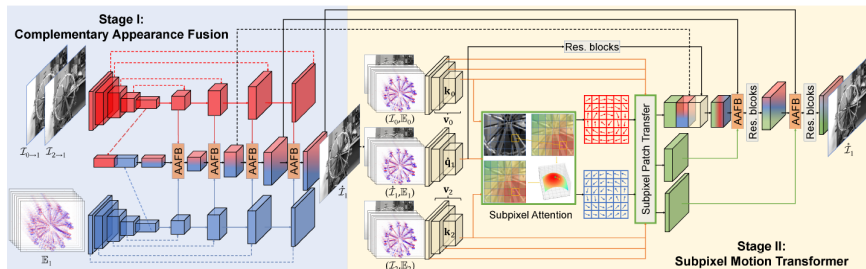


# Contributions

- Weakly supervised training, and perform SOTA, generalize well.
- Complementary appearance fusion(*CAF*) for aggregates image and event appearance.
- Subpixel Motion Transformer(*SMT*), for super resolution(*SR*) reconstruction.
- Contribute an event dataset, but not open source yet.

# Proposed Method

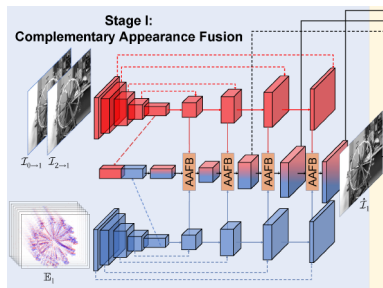
## Overview



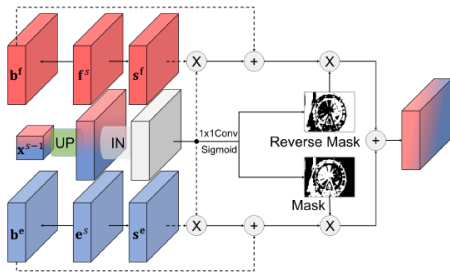
- Two stages: complementary appearance fusion(CAF) and subpixel motion transfer(SMT)
- CAF stage: fuse images and events, get initial  $\hat{I}_1$ .
- SMT stage: refine results with subpixel attention mechanism.

# Proposed Method

## Stage I: Complementary appearance fusion(CAF)



Stage I



AAFB

- Prepare works:

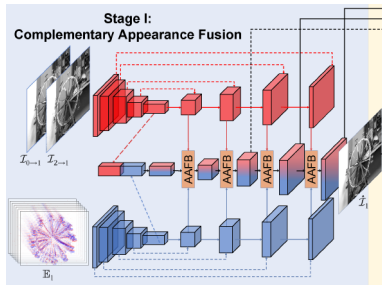
- ▶ Forward warping  $\mathcal{I}_0$  and  $\mathcal{I}_2$  with PWC-Net<sup>1</sup>  $\rightarrow \mathcal{I}_{0 \rightarrow 1}, \mathcal{I}_{2 \rightarrow 1}$
- ▶ Event representation  $\mathbb{E}_1$

<sup>1</sup>Deqing Sun et al. "Pwc-net: Cnns for optical flow using pyramid, warping, and cost volume". In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2018, pp. 8934–8943.

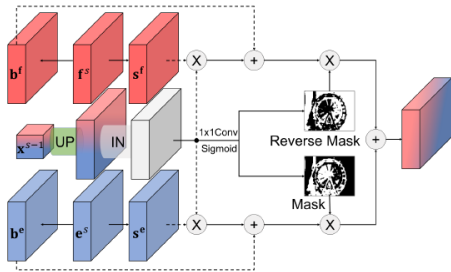


# Proposed Method

## Stage I: Complementary appearance fusion(CAF)



Stage I



AAFB

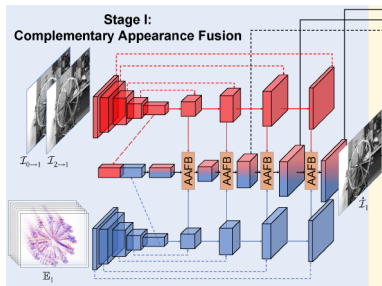
- Adaptive Appearance Fusion Block(AAFB):

$$x^s = g \left( x_{\uparrow}^{s-1}; f^s, e^s \right), s \in \{1, 2, 3, 4, 5\} \quad (1)$$

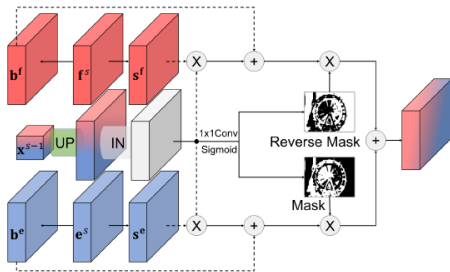
$x_{\uparrow}^{s-1}$  denotes the  $2 \times$  upsampled version of  $x^{s-1}$   
 $f^s, e^s$  denotes 2 branches features.

# Proposed Method

## Stage I: Complementary appearance fusion(CAF)



Stage I



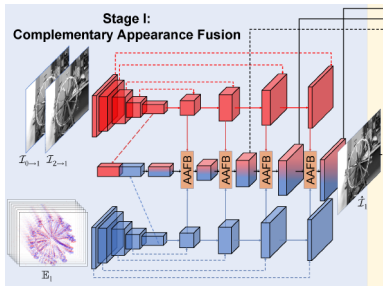
- Adaptive Appearance Fusion Block(AAFB):

$$y^e = \left( \frac{x_{\uparrow}^s - \mu(x_{\uparrow}^s)}{\sigma(x_{\uparrow}^s)} \right) \odot s^e + b^e \quad (1)$$

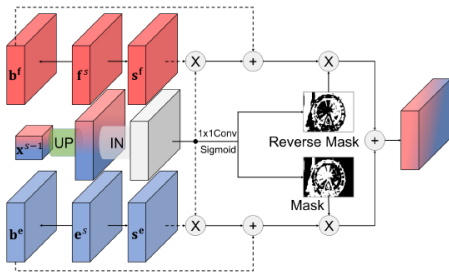
Break up  $f^s$ ,  $e^s$  with scalings and biases  $s^f$ ,  $b^f$  and  $s^e$ ,  $b^e$

# Proposed Method

## Stage I: Complementary appearance fusion(CAF)



Stage I



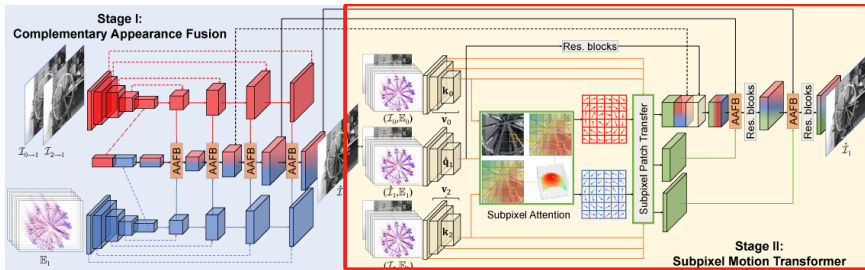
- Adaptive Appearance Fusion Block(AAFB):

$$y^e = \left( \frac{x_{\uparrow}^s - \mu(x_{\uparrow}^s)}{\sigma(x_{\uparrow}^s)} \right) \odot s^e + b^e \quad (1)$$

$$y = y^e \odot \mathbf{m} + y^f(1 - \mathbf{m}) \quad (2)$$

# Proposed Method

## Stage II: Subpixel Motion Transformer(SMT)



- Prepare works:

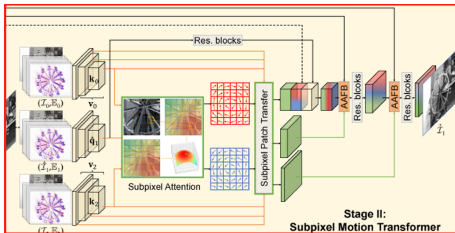
- ▶  $(I, \mathbb{E}) \xrightarrow{\text{conv}} \{v^s \mid s \in \{0, 1, 2\}\}$  as **values**
- ▶  $v^2 \xrightarrow{\text{clone}} k_0, k_2$  as **keys**,  $\hat{k}_1$  as **query**
- ▶ relevance measure:

$$D_0(i, p) = \left\| \frac{\hat{k}_1(i)}{\|\hat{k}_1(i)\|_2} - \frac{k_0(i + p)}{\|k_0(i + p)\|_2} \right\|_2^2 \quad (1)$$

$p \in [-m, m]^2$  represents a spatial **offset**.

# Proposed Method

## Stage II: Subpixel Motion Transformer(SMT)



- Subpixel attention learning:

$$j = i + p^* \text{ where } p^* = \arg \min_p D_0(i, p) \quad (1)$$

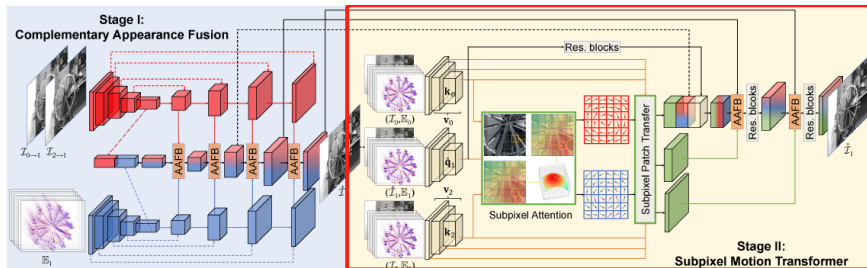
$$d(u) = D_0(i, p^* + u), u \in \mathbb{Z}^2 \cap [-n, n]^2 \quad (2)$$

$$d(u) \approx \hat{d}(u) = \frac{1}{2} u^T A u + b^T u + c \quad (3)$$

$$\min_{A, b, c} \sum_u w(u) \|\hat{d}(u) - d(u)\|^2 \quad (4)$$

# Proposed Method

## Stage II: Subpixel Motion Transformer(SMT)



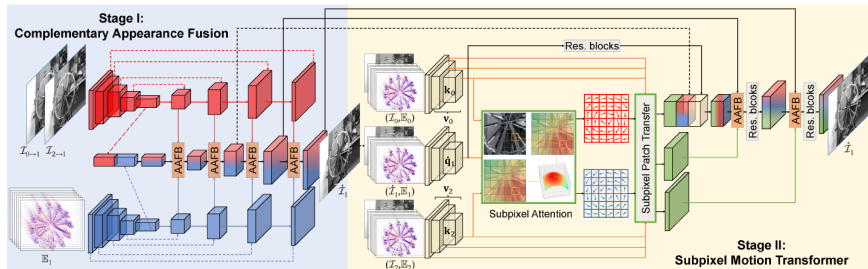
- Subpixel patch transfer:

- ▶ distance means relevance( $K \cdot Q$ ), then multiply with  $V$ , get transferred values  $z_0^s$  and  $z_2^s$
- ▶ select transferred values

$$z_1^s(i) = \begin{cases} z_0^s(i), & \text{if } D_0(i, p_0^*) < D_2(i, p_2^*) \\ z_2^s(i), & \text{otherwise.} \end{cases} \quad (1)$$

- ▶ fuse with previous features

# Proposed Method



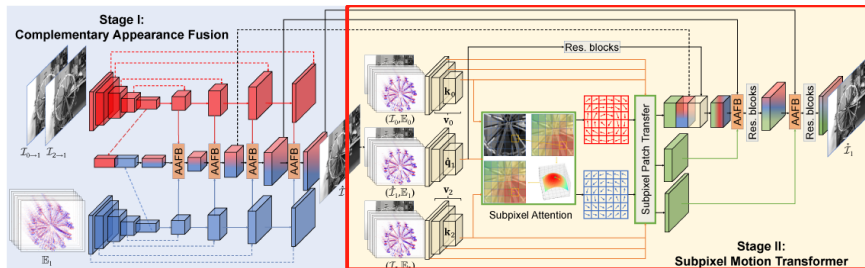
## Settings:

- ▶ Datasets: GroPro(720p, 240FPS), SloMo-DVS dataset(self-collected)  
Adopt  $10\times$ , so sample 1th, 11th, 21th frames to form a **sparse** training triplet to train.
- ▶ Loss: Charbonnier Loss(superior to  $\ell_1$ ,  $\ell_2$  loss)

$$\text{Charbonnier\_loss} = \sqrt{x^2 + \epsilon^2} \quad (1)$$

# Proposed Method

## Thoughts



- Frame info still rely on optical flow, not end-to-end.
- Process event in image format, point cloud or 3D conv works?
- 240FPS with sample rate 10  $\xleftrightarrow{\text{versus}}$  25FPS with sample rate 1, which is weakly supervision?



# Experiments

## Quantitative results

Table 1. Comparing models on GoPro dataset, measured in PSNR and SSIM. Bold indicates the top place while underline the second.

Supervision Methods	High FPS videos					High FPS videos + events					Low FPS videos + events	
	SloMo <sup>[14]</sup>	QVI <sup>[48]</sup>	DAIN <sup>[3]</sup>	TAMI <sup>†[7]</sup>	FLAVR <sup>[16]</sup>	ETV <sup>[38]</sup>	SloMo <sup>*[14]</sup>	QVI <sup>*[48]</sup>	EMD <sup>[15]</sup>	EDVI <sup>[21]</sup>	BHA <sup>[33]</sup>	Proposed
PSNR	27.79	29.54	27.30	32.91	31.10	32.25	32.79	<u>33.07</u>	29.67	30.90	28.49	<b>33.33</b>
SSIM	0.838	0.872	0.836	<b>0.943</b>	0.917	0.925	<u>0.940</u>	<b>0.943</b>	0.927	0.905	0.920	<u>0.940</u>

† TAMI also adopts external private datasets for training. \* Enhanced variants with events added into the inputs of network.

Table 2. Comparing models on SloMo-DVS dataset, measured in PSNR and SSIM. Bold indicates the top place while underline the second.

Supervision Methods	High FPS videos				High FPS videos + events				Low FPS videos + events	
	SloMo <sup>[14]</sup>	QVI <sup>[48]</sup>	DAIN <sup>[3]</sup>	FLAVR <sup>[16]</sup>	ETV <sup>[38]</sup>	SloMo <sup>*[14]</sup>	QVI <sup>*[48]</sup>	EDVI <sup>[21]</sup>	BHA <sup>[33]</sup>	Proposed
PSNR	30.69	30.93	30.38	30.79	32.06	33.46	<u>33.70</u>	33.60	22.95	<b>34.17</b>
SSIM	0.915	0.920	0.914	0.917	0.936	0.950	<b>0.953</b>	0.948	0.828	<u>0.952</u>

# Experiments

## Abalation Study

Table 3. Performance in PSNR with low frame-rate training.

Method	SloMo <sup>*[14]</sup>		QVI <sup>*[48]</sup>		Proposed
	High	Low	High	Low	
GoPro	32.79	31.40	33.07	29.88	33.33
SloMo-DVS	33.46	32.76	33.70	31.80	34.17

Table 4. Analysing the performance of CAF network.

Setting	PSNR	SSIM
Replacing AFFB with cat.+conv.	32.27	0.930
Using image branch only	29.43	0.882
Using event branch only	31.37	0.927
Full model	32.47	0.929

Table 5. Analysing the performance of SMT network.

ID	Key type	Value type	Att. type	Fused stage	PSNR
1	img.+evt.	image	subpix.	both	32.72
2	img.+evt.	event	subpix.	both	32.91
3	image	img.+evt.	subpix.	both	33.01
4	event	img.+evt.	subpix.	both	33.03
5	img.+evt.	img.+evt.	subpix.	first	33.00
6	img.+evt.	img.+evt.	subpix.	second	32.56
7	img.+evt.	img.+evt.	hard	both	33.02
8	img.+evt.	img.+evt.	soft	both	32.50
9	img.+evt.	img.+evt.	subpix.	both	33.33

# Experiments

## Qualitative results

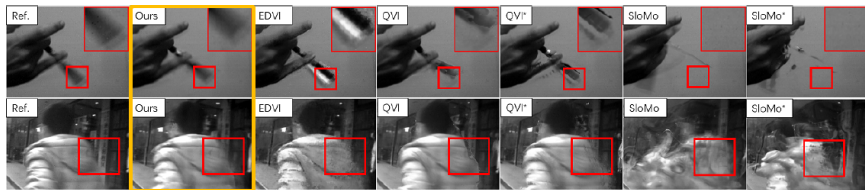


Figure 5. Qualitative comparisons on real data. In the first column (Ref.) we visualize the nearest input frame as reference since there is no groundtruth. We suggest the readers to watch our supplementary video for more qualitative comparisons on real-world video interpolation.

# Experiments

## Qualitative results

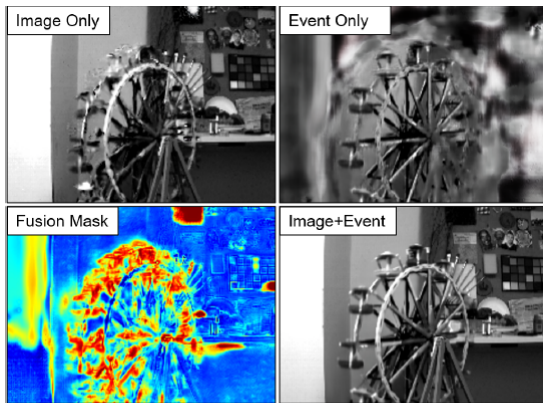


Figure 6. Visualizing the impact of adaptively fusing image and event appearance features in the CAF network.

# Experiments

## Qualitative results

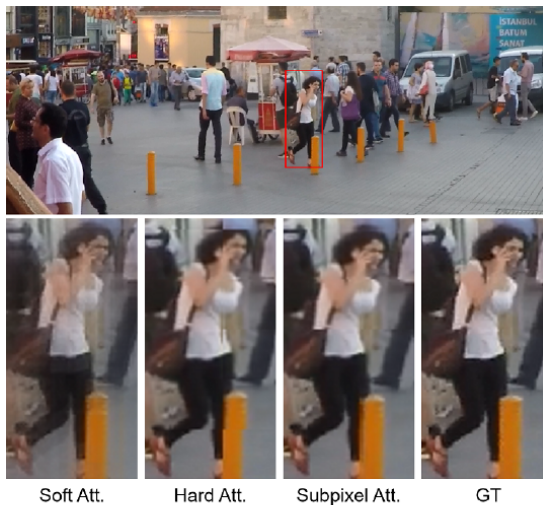
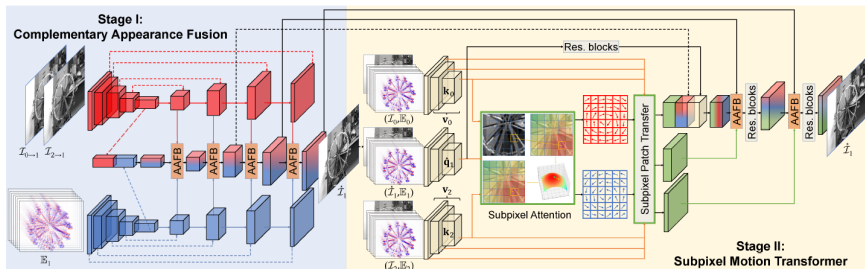


Figure 7. Patch transfer results with different types of attention.

# Conclusions



- Train with low frame-rate videos, but outperform others and generalize well.
- Attention mechanism, no direct motion modelling, eg optical flow.

*Thanks, Q & A*