Event-driven Video Frame Synthesis

Zihao Wang¹, Weixin Jiang¹, Kuan He¹, Boxin Shi², Aggelos Katsaggelos¹, Oliver Cossairt¹

¹Northwestern University

² Peking University







2nd Int'l Workshop on Physics Based Vision meets Deep Learning (PBDL) in conjunction with







Physics-based optimization

Subject to noise, or incomplete modeling



Physics-based vision







Physics-based optimization

Subject to noise, or incomplete modeling

 \oplus^*



Learning-based vision





End-to-end non-linear fitting

usually have superior performance as long as you have sufficient data and GPUs



Learning-based vision





End-to-end non-linear fitting

Learning from simulation!





Learning-based vision





End-to-end non-linear fitting

Train with data augmentation



Learning-based vision





End-to-end non-linear fitting

Require retraining even for similar tasks.







Frame-based camera pipeline



- We need "smart" cameras that:
 - Can respond to high speed motions (eliminate blur)
 - Do not always operate at high speed (less data redundancy)
- Potential solution: event cameras

What's event camera? Another high-speed camera?



Capture: 22 FPS

Display: 1.1 FPS

But...

- Events ≠ temporal gradient
 - Infinitely many solutions to infer intensity from events.
 - Cannot capture weak variations



But...

- Events ≠ temporal gradient
 - Infinitely many solutions to infer intensity from events.
 - Cannot capture weak variations
- Events are very noisy
 - Noise model not well understood
 - Gaussian on threshold
 - Event denoisers not advanced
 - Able to cancel isolated events (correlation)
 - Cannot handle complex scenarios, e.g. illumination change

Example images overlaid with neighbor events data from DAVIS dataset and Pan et al. CVPR'19



We propose intensity frame + events for high frame-rate video synthesis



Our approach: fusion of intensity frame + events



Differentiable model (event sensing)



Differentiable model (frame sensing)

We consider 3 temporal settings: interpolation, prediction and motion deblur.

Case (intensity tensor notation)	Model
Interpolation (\mathcal{F}^i)	$\mathscr{F}_1^i=\mathscr{H}_1, \mathscr{F}_2^i=\mathscr{H}_d$
Prediction (\mathscr{F}^p)	$\mathscr{F}^p=\mathscr{H}_1$
Motion deblur (\mathscr{F}^m)	$\mathscr{F}^m = rac{1}{d} \sum_{t=1}^d \mathscr{H}_t$



Reconstruction loss and optimization

Objective
$$\hat{\mathscr{H}} = \operatorname*{argmin}_{\mathscr{H}} \mathcal{L}_{pix}(\mathscr{H}, \mathcal{F}, \mathcal{E}) + \mathcal{L}_{TV}(\mathscr{H})$$

Frame pixel error Pixel loss $\mathcal{L}_{pix}(\mathscr{H}, \mathcal{F}, \mathcal{E}) = \mathbb{E}_{fpix}[\|\mathcal{F} - \mathcal{A}(\mathscr{H})\|_1]$ $+ \lambda_e \mathbb{E}_{epix}[\|\mathcal{E} - \mathcal{B}(\mathscr{H})\|_1]$ Event pixel error

Sparsity loss
$$\mathcal{L}_{TV}(\mathscr{H}) = \lambda_{xy} \mathbb{E}_{hpix} \left[\left\| \dot{\mathscr{H}}_{xy} \right\|_{1} \right] + \lambda_{t} \mathbb{E}_{hpix} \left[\left\| \dot{\mathscr{H}}_{t} \right\|_{1} \right]$$

 $\dot{\mathscr{H}}_{xy} = \frac{\partial \mathscr{H}}{\partial x} + \frac{\partial \mathscr{H}}{\partial y} \quad \dot{\mathscr{H}}_{t} = \frac{\partial \mathscr{H}}{\partial t}$

Use stochastic gradient descent (SGD) to minimize the loss. As loss decreases, results get closer to ground truth.



Results (DMR)

- Interpolation case
 - Given start & end frames + events in-between, recover intermediate frames

Low-speed intensity frames (2 frames)

Event frames (20 frames)

High-speed video (21 frames)



The middle frame is withheld for evaluation







Results (DMR)

- Prediction case
 - Given start frame and future events, recover future frames



CF [ACCV'18]

PSNR: 23.33 SSIM: .771





PSNR: 25.12 SSIM: .831



PSNR: 36.59 SSIM: .983





Results (DMR)

- Motion deblur case
 - Given a blurry image + events in-exposure, recover intermediate sharp frames.

Blurry images



Events during exposure





EDI [CVPR'19]



Ours

Video recovery





Overview of our approach



Residual "denoiser"

- Use CNN to learn the residual of DMR output w.r.t. ground truth
 - Designed to enhance DMR results
- Easy to train
 - Model DMR artifacts as residual "noise"
 - Actually beyond Gaussian denoising
- Single frame based
 - Interface well with DMR



Results (residual denoiser)





clip name	metric	DMR	DnCNN	FFDNet	Ours
airplane	PSNR	30.91	31.10	30.92	31.38
	SSIM	0.975	0.982	0.976	0.982
basketball	PSNR	23.55	24.05	23.47	24.06
	SSIM	0.963	0.971	0.964	0.972
soccer	PSNR	29.96	31.08	30.13	31.29
	SSIM	0.961	0.974	0.962	0.975
billiard	PSNR	36.46	35.42	36.48	36.46
	SSIM	0.982	0.986	0.983	0.987
ping pong	PSNR	32.46	32.26	32.50	32.24
	SSIM	0.974	0.978	0.975	0.979

10/29/2019

Results

- Comparison with non-event-based frame interpolation approach
 - Events can provide additional information which is useful for challenging motions.

SepConv [CVPR'17]



Ground truth



Ours (DMR + RD)



