

Event camera aided video frame interpolation

- a report for CS396 Natural and Artificial Vision 2023 Spring

Daizong Tian^a

This manuscript was compiled on June 4, 2023

Conventional camera based video frame interpolation (temporal video super resolution, or VFI in this report) has been a very popular topic in Computer Vision area. For optical flow based VFI methods, one of the major challenge is to compensate for inaccurate flow generated from optical flow estimation modules (mostly neural networks). Optical flow estimation is an ill-posed problem since it's only estimated based on pixel values and a pre assumption need to be made on object motions. Event camera can helps VFI since it can record a much continuous motion in terms of high frame-rate events thus could provide useful information for VFI and brings little extra costs.

Event camera | Video frame interpolation | Super resolution | Optical flow

Trade-offs between video resolution (spatial quality) and frame rate (temporal quality) has always been an issue for producers and consumers. The trade-off can be coming from the limitation of sensor read-out speed for video recorders, the computational power from gaming consoles, the bandwidth for online video streamers, etc... To solve this problem, people usually choose to capture/generate a high spatial but low temporal resolution video first and then generate more frames when displaying it out of the existing frames. This is actually also similar to how inter-frame video compression works: select some of the frames as high quality I frames and the rest frames (B and P) only contains information that changed compared the the reference I frames. In VFI, it's called MEMC: motion estimation and motion compensation, which is very popular in TVs. Recent years NN-based VFI becomes more and more popular. NVIDIA provide Super resolution (both spatial and temporal) support for their video cards.

However, conventional camera frame-based video frame interpolation have drawbacks as described in the abstract. Luckily event camera perfectly provide the information that conventional camera missed but critical to VFI: the motion information. Nowadays event-camera and conventional camera combined VFI methods can achieves far better results than conventional camera-only VFI sota methods. (1–3)

Event camera, a bio-inspired motion extractor

To people's surprise who firstly nows this area, event-camera are human-vision-sytem (HVS) inspired. Event camrea are also called silicon retinas sometimes and it actually mimics some part of our HVS. Our human retina are sensitive to light intensity changes instead of absolute light intensity, which can be side proved by the facts that there are blood vessels and neurons in front of our retina but we cannot see them: Since the tissues relative position to the retina does not change, it does not provide any light intensity changes so our retinas ignores it.

We all familiar with conventional cameras, that captures photo-generated electron within a certain amount of time to generate voltage, then convert it into digital pixel values. In other words, get the absolute light intensity within a certain amount of time. The event camera, however, does not capture the absolution intensity, it only captures if the light intensity becomes bright or darker to

a certain step, and outcome with on or off events for every pixel location (and no event if light intensity not changing). This brings two benefits, 1: it does not need to pre-set a shutter speed so events generates very fast, since a moving object are usually continuously brings light intensity changes. And 2: A large part of the scene usually does not change, so event camera will not and cannot captures them, that saves data bandwidth and computational power. Those two special features makes event camera very high frame rate (or event rate) while maintains low power-consumption (usually 1/1000 of conventional camera sensors) (4).

One can imagine a event camera is a high-passed version of conventional high speed camera with a even higher speed of capturing the motion. Although reconstruct scene directly from event camera is also a hot topic in Computer vision and computational imaging, current combining conventional camera and event camera together usually brings a better results since we can use the good pixel information from conventional camera and good motion information from event camera.

Overview of frame based, optical flow based video frame interpolation

Modern optical flow based conventional camera video frame interpolation can mostly divided into two parts: optical flow estimation (motion estimation) and frame interpolation. FlowNet (5), pwc-net (6), and RAFT (7) are some popular optical flow estimation networks. Most networks estimate optical flow by two input frames based on pixel values or first or second order statistics, which is a ill-posed problem since the time gap between two frame are two much so the possible motion field is very large, not to mention differencnt object may share same statistics and same object may have different statistics in different locations. Thus conventional camera based optical flow estimation gives inaccurate flows especially on smooth areas like backgrounds.

To overcome this problem, video frame interpolation methods that based on optical flow estimation added additional modules in interpolation part to correct and compensate the inaccurate optical flows. As an example, Super-Slomo (8) not only have a optical residual module to correct inaccurate optical flows generated from FlowNet (5) (and error introduced during flow-reversal for backward warping), two of the four loss components are used to regulate flow estimation.

Another problem is that due to the huge gap between frames, the actual motion can never be re-constructed but only guessed based on assumptions. Super-Slomo (8) assumes linear motion so intermediate frame will generate object in between of locations of the previous and following frames. Quadratic Frame Interpolation (QVI) (9) assumes quadratic motion with the fact that a free-fall object will receives gravity so will have a quadratic motion, and thus needs 4 frames to estimate this quadratic motion with

Author affiliations: ^aNorthwestern Univerisity

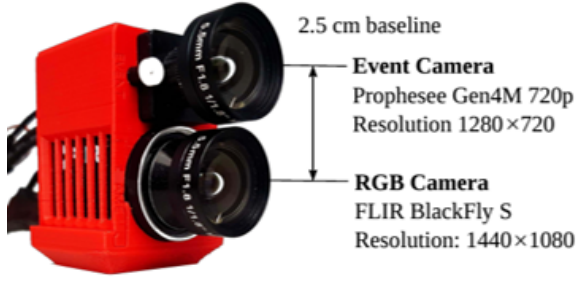


Fig. 1. Structure of TimeLens (1) proposed camera system, constructed with a conventional RGB camera and an event camera

the assumption that the object motion among the 4 frames are continuous. However, those assumptions will not always agree with the ground-truth motion and it cannot possibly be truly reconstructed unless additional motion information is given.

Event camera aided Video Frame Interpolation

The event camera can fill in the gap that conventional camera-only video frame interpolation lacks. The event camera is a high-pass filtered camera that contains and only contains motion information, which is exactly the optical-flow-based VFI needed. In addition, the events already contain the pixel change information, thus extracting motion features (optical flow) by a NN and manually warping pixels might be unnecessary, synthesis-based video frame interpolation can simply take pixel info (frames) and pixel change info (events) to generate intermediate frames.

Optical-flow-based interpolation, synthesis-based interpolation, both have their pros and cons.

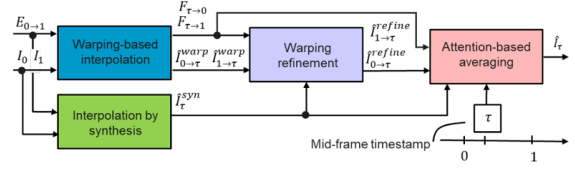
Optical flow based interpolation, as we mentioned earlier, generates a dense representation of motions and then warps pixel values from the original to destinations. It assumes pixel values of the same object not changing, thus it is struggling to handle illumination change and object reshaping. However, since it's based on dense motion maps and strict warping, the output quality could be ensured if the optical flow map is accurate, just might with wrong illumination.

Synthesis-based interpolation, on the other hand, does not make any assumption. Thus it can handle illumination and shape change if considering an ideal neural network. However, in the real world, due to lacking of constraints compared to optical-flow-based methods, it sometimes generates undesired artifacts.

Fig. 3 shows a comparison between two methods and the ground truth image. TimeLens (1) combines the good from both methods, and proposed a conventional camera and event camera bounded video frame interpolation techniques. So I will use it as an example of how event cameras can help VFI. The hardware structure is shown in Fig. 1.

And it uses several Neural Networks with designed structures to combine the benefits from both conventional and event cameras, both optical-flow-based interpolation and synthesis-based interpolation together. The overall design is shown in Fig. 2 The backbone of all NNs are modified hourglass networks.

The warping-based module will use an hourglass network to generate optical flow using event-camera-generated data, warping the previous and following frame to the desired time spot to generate two separate frames. The synthesis-based network will use another hourglass network to receive both event-camera data and conventional camera RGB frames to generate a frame



(a) Overview of the proposed method.

Fig. 2. Structure of TimeLens (1), every colored box is a module with a separate Neural Network

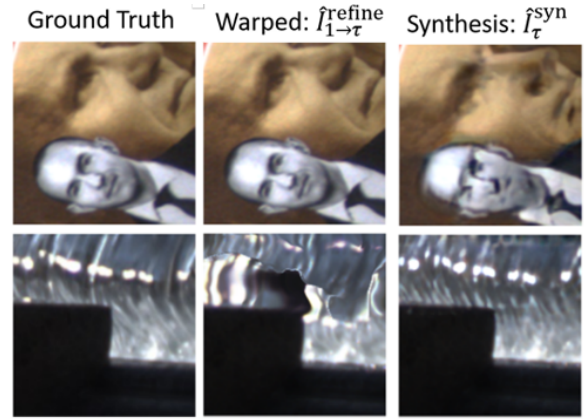


Fig. 3. Ground Truth V.S. Warping-based method V.S. Synthesis-based method. (1)

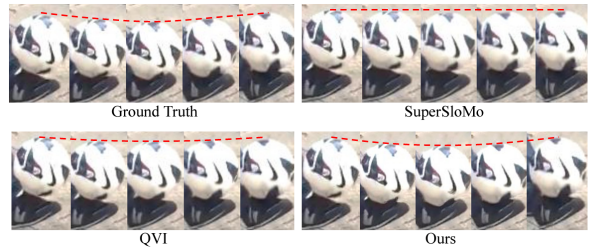


Fig. 4. Time Replayer (2) V.S. conventional camera VFI

at desired time spot. Then the synthesised result and warped result and optical flow will be fed into another hourglass network compensate possible optical flow errors like misalignment (notice that there are 2.5cm baseline in Fig.2 between two cameras). In the following attention module, there's a fourth hourglass network that takes both warped frames, optical flow, and synthesised frame to generate weights to average the 3 generated frames to get a final interpolated output. With 4 networks and information from both cameras, it achieves around 4-9 dB increment in terms of PSNR on popular datasets like GoPro(10), Vimeo90k(11). Notice that since those dataset does not provide event data, the events are acutally simulated using an event simulator (12) with skipped frames.

TimeReplayer (2), another event and conventional camera combined VFI, proves that the combined method can record and thus interpolate the actual object movements. Fig.4 clearly shows the "Ours" (which means TimeReplayer's method(2)) recovers the football motion much better than linear-assumption SuperSloMo(8) and quadratic-assumption QVI(9).

Conclusion

The event camera, with its benefits of high temporal resolution and low data rate, perfectly fills the need for conventional camera video

frame interpolations. I truly believe it will make VFI much easier with less power consumption and higher quality, and can come into a real-world application like how ToF lens becomes standard in today's smartphones.

I have slides also publicly viewable on <https://makiseasuka.github.io/>

Ethics

As an image-based neural network, this area has some potential ethics concerns: Firstly, it raises privacy issues as the system may inadvertently capture and process images that individuals may not wish to be shared. Moreover, inherent biases within the data used to train these neural networks can lead to biased outputs, potentially causing inequities or unfair treatment, especially considering event cameras might be more sensitive to some colors since it only records intensity changes and different colors may have different scale factor in terms of colors.

Therefore strict privacy regulations should be implemented regarding data collection and usage. Although there are public datasets available, only verified datasets should be used, also, to address bias, it's important to ensure that training datasets are diverse and representative of the populations. After training, detailed tests should be conducted for all different situations.

- 1 S Tulyakov, et al., Time lens: Event-based video frame interpolation in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. pp. 16155–16164 (2021).
- 2 W He, et al., Timereplayer: Unlocking the potential of event cameras for video interpolation in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 17804–17813 (2022).
- 3 J Dong, K Ota, M Dong, Video frame interpolation: A comprehensive survey. *ACM Transactions on Multimed. Comput. Commun. Appl.* (2022).
- 4 G Gallego, et al., Event-based vision: A survey. *IEEE transactions on pattern analysis machine intelligence* **44**, 154–180 (2020).
- 5 E Ilg, et al., FlowNet 2.0: Evolution of optical flow estimation with deep networks in *Proceedings of the IEEE conference on computer vision and pattern recognition*. pp. 2462–2470 (2017).
- 6 D Sun, X Yang, MY Liu, J Kautz, Pwc-net: Cnns for optical flow using pyramid, warping, and cost volume in *Proceedings of the IEEE conference on computer vision and pattern recognition*. pp. 8934–8943 (2018).
- 7 Z Teed, J Deng, Raft: Recurrent all-pairs field transforms for optical flow in *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part II 16*. (Springer), pp. 402–419 (2020).
- 8 H Jiang, et al., Super slo-mo: High quality estimation of multiple intermediate frames for video interpolation in *Proceedings of the IEEE conference on computer vision and pattern recognition*. pp. 9000–9008 (2018).
- 9 X Xu, L Siyao, W Sun, Q Yin, MH Yang, Quadratic video interpolation. *Adv. Neural Inf. Process. Syst.* **32** (2019).
- 10 S Nah, T Hyun Kim, K Mu Lee, Deep multi-scale convolutional neural network for dynamic scene deblurring in *Proceedings of the IEEE conference on computer vision and pattern recognition*. pp. 3883–3891 (2017).
- 11 T Xue, B Chen, J Wu, D Wei, WT Freeman, Video enhancement with task-oriented flow. *Int. J. Comput. Vis.* **127**, 1106–1125 (2019).
- 12 D Gehrig, M Gehrig, J Hidalgo-Carrió, D Scaramuzza, Video to events: Recycling video datasets for event cameras in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 3586–3595 (2020).