

Neural Ganglion Sensors: Learning Task-specific Event Cameras Inspired by the Neural Circuit of the Human Retina

Haley M. So and Gordon Wetzstein
Stanford University

Abstract

Inspired by the data-efficient spiking mechanism of neurons in the human eye, event cameras were created to achieve high temporal resolution with minimal power and bandwidth requirements by emitting asynchronous, per-pixel intensity changes rather than conventional fixed-frame rate images. Unlike retinal ganglion cells (RGCs) in the human eye, however, which integrate signals from multiple photoreceptors within a receptive field to extract spatio-temporal features, conventional event cameras do not leverage local spatial context when deciding which events to fire. Moreover, the eye contains around 20 different kinds of RGCs operating in parallel, each attuned to different features or conditions. Inspired by this biological design, we introduce Neural Ganglion Sensors, an extension of traditional event cameras that learns task-specific spatio-temporal retinal kernels (i.e., RGC “events”). We evaluate our design on two challenging tasks: video interpolation and optical flow. Our results demonstrate that our biologically inspired sensing improves performance relative to conventional event cameras while reducing overall event bandwidth. These findings highlight the promise of RGC-inspired event sensors for edge devices and other low-power, real-time applications requiring efficient, high-resolution visual streams.

1. Introduction

Event cameras can offer numerous advantages over frame-based image sensors that are crucial for the extreme constraints of emerging edge devices such as autonomous vehicles, robotics, and augmented/virtual reality. These include high temporal resolution, low latency, low power consumption, low bandwidth, and high dynamic range. Whereas traditional cameras output intensity frames at a fixed frame rate, a traditional event camera outputs asynchronous spikes that capture *differences* in intensity. This design was initially inspired by the spiking nature of retinal ganglion cells (RGCs) which encodes the light hitting our retina into information-dense spikes and transmits the signals to our brain. Event cameras have shown promise in many tasks including action recognition [21, 22, 45], optical flow

[51, 52, 58], depth or shape estimation [24, 39, 60], and more [20]. However, event cameras are still limited compared to their biological analogues. While RGCs can aggregate information spatially across an area on the retina, event cameras only operate on a per-pixel basis: each pixel sends events independent of neighboring pixel information.

Most RGCs use local spatial information to decide whether or not to fire. There are roughly 20 kinds of RGCs in the human eye, and while the exact function of each varies, in general, these retinal ganglion cells receive information from not just a single photoreceptor, but rather a small group (as in midget cells) or from an even larger receptive field (as in parasol cells) to decide whether to send an action potential to the brain [61]. One identified organization these RGCs operate with is a center-surround organization, in which the center receptive field is compared to the surrounding larger selected region to determine whether or not to fire a spike [29]. For example, CENTER ON cells look for where the center is on (light falls on the photoreceptor) and the surround is off, and CENTER OFF cells look for the inverse. In the human eye, this organization can allow for edge and contrast enhancement. In addition, there are direction-selective RGCs that are attuned to specific spatial frequencies, color-sensitive RGCs, and temporally attuned RGCs specialized for flicker or motion, etc. [50]. There are roughly 1.5 million RGCs in our eye that distill the information falling on our roughly 120 million photoreceptors, encoding the spatio-temporal information into sparse binary spikes. This raises our motivating question: How can we replicate the retina’s diverse retinal circuitry to achieve more efficient and versatile vision and perception?

Towards this end, we revisit the retina’s bandwidth-efficient design principles and seek to learn the optimal set of functions for generating events, bridging the gap between conventional event cameras and their biological counterparts. Specifically,

- we introduce a framework for learning asynchronous, spatio-temporal events to optimize performance and bandwidth for perception tasks;
- we develop a differentiable event simulator;
- we show that integrating local spatial information can im-

prove the performance on vision tasks over event cameras while exhibiting lower bandwidth;

- we explore learning multiple complementary event “channels,” further enriching the captured information and boosting performance.

2. Related Work

2.1. Retinal ganglion cells

Recent literature has made significant progress in understanding how visual signals are processed even before they leave the human retina [19, 25]. Here, photoreceptors sense the incoming light, and the resulting signals are transmitted through and modulated by a diversity of horizontal, bipolar, and amacrine interneurons before being processed by roughly 20 different types of retinal ganglion cells (RGCs). The variation in interneuron pathways and RGC activation create multiple parallel retinal circuits or functions that extract a variety of complementary visual features at the retina [49], resulting in spikes sent to the brain. While the specific visual features identified by each RGC can vary, the output of RGCs can generally be described by similar structures. As a result, biologists often use simplified models such as the linear–nonlinear (LN) cascade model [9, 13], the Hodgkin–Huxley model [26], and the various integrate-and-fire models [28, 44, 59]. There are also recent works in using deep learning techniques to try to predict retinal responses to natural images [7, 38, 43]. From this literature, the key insight that we apply to our work is that ganglions receive signals not only across time but also across space from a receptive field of photoreceptors, not just a single one. In addition, our eyes also use multiple types of RGCs in parallel for efficient sensing.

2.2. Event cameras

Mahowald and Mead first introduced the silicon retina in the late 1980s [35], which quickly led to the development of the Dynamic Vision Sensor (DVS) or the event camera [16, 32, 47]. In these systems, each pixel operates independently, always comparing the current intensity to a previously memorized intensity. If the intensity difference is larger than a set threshold value, an “event” is sent. The outputted information includes an x-y location, timestamp, and the polarity of the brightness change. The pixel then updates its memorized value to the intensity that triggered the output. Other variants of event cameras include the Asynchronous Time-based Image Sensor (ATIS) [46], where an event trigger will also readout the intensity value at the given pixel or encode intensity through inter-spike time intervals [14], and the Dynamic and Active Pixel Vision Sensor (DAVIS) [3, 10], which outputs full frame intensity values at a slow frame-rate along with the asynchronous events. While the events in these systems are inspired by

the basic spiking nature in the retina, they don’t utilize the spatial surroundings. Recently, Li and Delbrück introduced the Center Surround Dynamic Vision Sensor (CSDVS) and illustrated the potential benefits of using a center-surround organization and suggested a future hardware implementation to achieve such a sensor with lateral polysilicon resistors and controllable transverse conductance [17, 31]. In 2023, [33] proposed using a pattern of different event thresholds across the sensor, analogous to spatially varying pixel exposures seen in computational imaging [37, 41, 42]. Most recently in 2024, [57] presented Generalized Event Cameras, which explored a few kinds of statistical “differentiating” methods, of which included the spatial dimension. They also introduced a nice breakdown of a general event camera as “when to send” and “what to send.” While there are similar notes, the fundamental formulation as well as outputs are different: they output full intensity values while our approach remains truer to the human retina and the original event camera, sending just binary bits. Inspired by the human retina, we also uniquely explore learning multiple parallel task-specific RGCs and optimize specifically for bandwidth, not just performance.

2.3. Event-based processing in vision applications

There are a number of representations used to process asynchronous events in computer vision including event-by-event processing, events paired with intensity frames, and events processed in voxel grids. Recent review papers [20, 66] offer good insight into the advantages and disadvantages of each. Of these approaches, two of the most common ways to process events are either by utilizing spiking neural networks (SNNs) or by aggregating events into image-like frames and using conventional convolution neural networks (CNNs) [15]. While SNNs have been successful in a number of tasks including image classification [68], object detection [55], video reconstruction [70] and more [1, 8, 23, 67], they can be difficult to train as there is no traditional back propagation and they are relatively new, so are still catching up to the more mature field of CNN-based computer vision. As a result, most state-of-the-art networks [18, 36, 48, 67, 69] opt to aggregate events into voxel-like grids and use more traditional image-based computer vision techniques. As this is the most popular method, we create a differentiable binning procedure to backpropagate through event voxel grids to be able to learn task-specific events. In our work, we demonstrate that augmenting event cameras with additional biologically inspired spatial aggregation can improve performance for machine vision tasks. Furthermore, we explore learning multiple RGC channels and introduce a framework for optimizing these asynchronous events for bandwidth and performance.

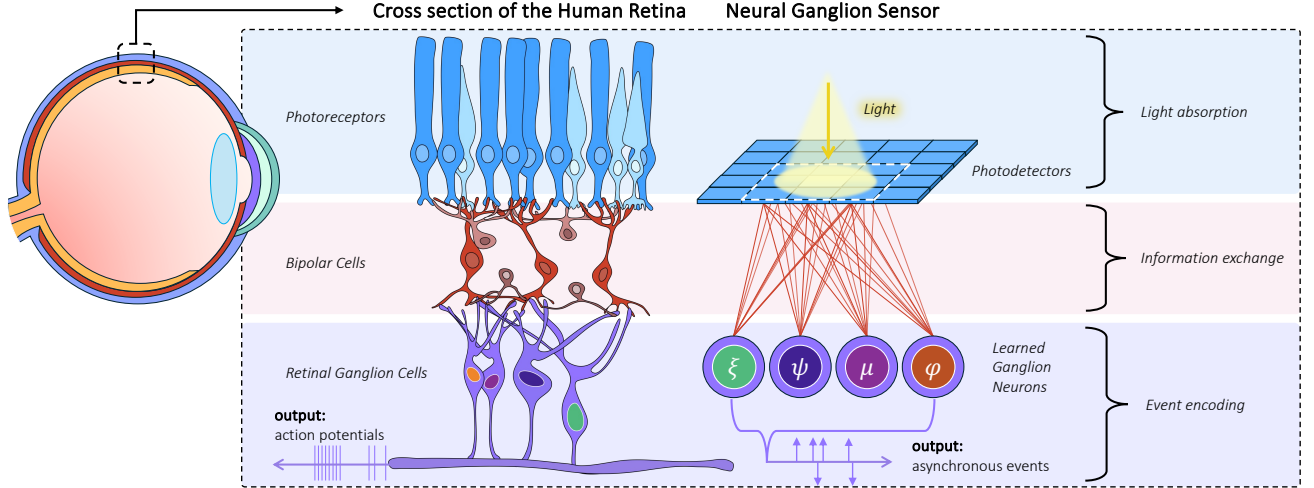


Figure 1. **Analogy between Neural Ganglion Sensors and the Human Retina.** On the left, we show a simplified diagram of different layers in the human retina. Light hits the photoreceptors (rods and cones), of which there are about 100 million per eye. The signals get transferred and modulated through Bipolar cells along with additional Horizontal and Amacrine cells. In the end, the roughly 1 million Retinal Ganglion Cells (RGCs), receive signals from a small *area* on the retina, not just from a single photoreceptor. These RGCs look at the pattern of information to decide whether to send a spike signal to the brain. We see spatial and temporal pooling occurs in the first few layers of the retina to encode all the information into bandwidth efficient spiking potentials. On the right, we show our proposed Neural Ganglion Sensor, an event camera augmented to better match the human retina.

3. Proposed Method

3.1. Model of RGC Events

The linear-nonlinear cascade (LN) is a popular way to model the response of RGCs [13]. In these models, the Poisson spike rate of a neuron is determined by a linear spatio-temporal filter and a nonlinear activation. Specifically, the probability, P that a neuron spikes can be described as a continuous 3D convolution of the intensity, I , over space x, y and time t followed by a non-linearity f :

$$P(x, y, t) = f\left(\left[W * I\right]_{(x,y,t)}\right) \quad (1)$$

where $*$ denotes the convolution operation, and W is a kernel with the weights of the spatiotemporal filter.

From the LN model, the event camera emerged with a few simplifications. Firstly, the probabilistic spiking is replaced with a deterministic activation if the output of f exceeds a threshold δ . Secondly, the continuous intensity I is split temporally into the current intensity I_{curr} and a memory intensity I_{mem} . Lastly, the W filter is reduced to a simple temporal differencing. For Neural Ganglion Sensors, our RGC event formulation uses the first two simplifications of conventional event cameras, but reintroduces the spatial dimension. The resulting model for when a RGC event is triggered becomes:

$$P(x, y, t) = \mathbb{1}\left(f\left(\left[W * (I_{curr} - I_{mem})\right]_{(x,y,t)}\right) > \delta\right) \quad (2)$$

where $\mathbb{1}$ is the indicator function.

Biological neurons are limited to outputting binary spikes, which is bandwidth efficient and fast to transmit. Similarly, when an event is triggered in event cameras, the output O for each event triggered is binary and can be computed by the polarity of the difference:

$$O(x, y, t) = \begin{cases} \text{sign}(I_{curr} - I_{mem})_{(x,y,t)} & \text{if } P(x, y, t) = 1 \\ \text{None} & \text{if } P(x, y, t) = 0 \end{cases} \quad (3)$$

In addition, when the event is triggered at pixel (x, y) , $I_{mem}(x, y)$ is updated to $I_{curr}(x, y)$.

$$I_{mem}(x, y) = \begin{cases} I_{curr}(x, y) & \text{if } P(x, y, t) = 1 \\ I_{mem}(x, y) & \text{if } P(x, y, t) = 0 \end{cases} \quad (4)$$

With these three equations defining the event trigger, the output, and the memory update, we have our full RGC event model. In fig. 1, we draw the parallel between the human retina and Neural Ganglion Sensors.

3.1.1. Event Camera model

Pixels in traditional event cameras operate independently of other pixels. Our RGC event formulation in equation 2 reduces to the traditional event camera when W is simply the identity kernel, and f is the absolute value function.

$$P(x, y, t) = \mathbb{1}\left(|I_{curr} - I_{mem}|_{(x,y)} \geq \delta\right) \quad (5)$$

The equations for the output (eqn. 3) and the memory update (eqn. 4) remain the same. While this formulation enables event cameras to mirror the temporal aggregation of the retina, ganglion cells observe light from a receptive field, not just at a single point [49].

3.1.2. Center-Surround Model

One type of spatial aggregation for RGCs is the center-surround organization. At a high level, this allows for contrast or edge enhancements, spatial filtering, and more. In a given receptive field, the center pixels are compared to the surrounding pixels. Center ON cells look for when the center is excited and the surround is inhibited. Center OFF cells look for the inverse. Midget and parasol cells, two of the most ubiquitous cell types in the retina, each have Center ON and OFF configurations, though their spatial reaches differ, making them more attuned to different spatial frequencies. To achieve this spatial behavior with our RGC event formulation, the center pixel can be compared to the average value of the surround. W would be set to a $k \times k$ kernel with 1 in the center pixel and weights summing to -1 in the surrounding pixels to model a Center ON. Center OFF would be the same kernel, but negated. Again, f is the absolute value function. Similarly, another variation of the center-surround was suggested by [17]. In their case, the W kernel was set to:

$$W = \begin{bmatrix} 0 & 1 & 0 \\ 1 & -4 & 1 \\ 0 & 1 & 0 \end{bmatrix} \quad (6)$$

3.1.3. Modeling other Human RGCs

With our RGC model, we can represent a number of other diverse RGC types found in the human eye. Although the exact functions of all human RGC types are still being characterized, incorporating spatial context plays an integral role in event sparsification. Our approach is not constrained to strictly biologically accurate RGCs but instead draws conceptual inspiration from their functional diversity.

3.2. Simulator to Learn Task-Specific RGC events

While the human retina inspires us, in the end, what visual features our eyes evolved to be attuned to may be very different from what machine vision would find important. As in deep learning where we moved from handcrafted features to learned features, here we follow a similar trajectory and learn what may be useful for perception tasks. We take our proposed model of a RGC event and learn the W kernel and the threshold values to tailor our sensing for vision applications. In order to do so, we need to backpropagate through to the event generation, so we build a differentiable event simulator which will be open-sourced.

3.2.1. Simulating RGC events

We construct our event generation based on ESIM and V2E [27], the two most widely used event simulators. However, we also support spatial kernels, allowing us to simulate

diverse types of RGC kernels that goes beyond conventional events. We incorporate non-uniform thresholds, refractory periods, separate and learnable positive and negative contrast thresholds, and the option to operate on log or linear intensity. Analogous to how human eyes have 20 kinds of RGCs, we also support learning multiple kinds of events, either through spatially-varying thresholds and kernels in a 2-by-2 bayer-like pattern or through multi-channel events.

3.2.2. Differentiable binning

Current state-of-the-art event-based vision pipelines pre-process events into sparse frame-like images or sparse voxels. This allows researchers to build off of the plethora of works in frame-based computer vision. In this work, we seek to learn events for two different tasks, one that uses events and RGB images, as is common with the DAVIS sensor and other emerging industry sensors, and the other that uses just events. In both cases, the state-of-the-art works we build off of pre-process the events into a voxel-like data structure. To learn our RGC kernels, we must backpropagate through these voxel structures to the event generation, so we implement differentiable binning. Our binning is performed with a closed-form solution that can compute binned RGC events from high-speed video. This combines the functionalities of video-to-event simulators and conventional event pre-processing. Following recent approaches [56, 63], events are weighted linearly into the two closest time bins. For each pair of frames at time t_k and t_{k+1} in the input high-speed video, the events generated between the frames are distributed into the two nearest neighboring bins at times t_{bin}^- and t_{bin}^+ , as described by the following equations. Here, I_{RGC} is the output of the RGC kernels, α is the time spacing between events, β is the time offset to bin^- , pol is the polarity of the generated events, and \mathcal{N} is the quantized number of generated events:

$$\begin{aligned} I_{RGC} &= W * (I_{curr} - I_{mem}) \\ pol &= \text{sign}(I_{RGC}) \\ \alpha &= \frac{t_{k+1} - t_k}{I_{RGC}/\delta} \\ \beta &= (t_k - t_{bin}^-) \\ \mathcal{N} &= I_{RGC}/\delta \end{aligned} \quad (7)$$

where δ is either the positive or negative threshold, depending on the polarity of I_{RGC} . Each event generated is weighted linearly into the two bins, depending on the distance, w_i , to each bin.

$$\begin{aligned} w_i &= \frac{\beta + (i + 1) * \alpha}{t_{bin}^+ - t_{bin}^-} \\ bin_i^- &= (1 - w_i) \cdot pol \\ bin_i^+ &= w_i \cdot pol \end{aligned} \quad (8)$$

Binning all the events generated between the pair of frames, we get the following:

$$\begin{aligned} \text{bin}_{frame}^- &= \sum_{i=0}^{\mathcal{N}-1} \text{bin}_i^- \\ \text{bin}_{frame}^+ &= \sum_{i=0}^{\mathcal{N}-1} \text{bin}_i^+ \end{aligned} \quad (9)$$

We can derive the closed-form solution for these summations by using the formula for the sum of an arithmetic sequence:

$$\begin{aligned} \text{bin}_{frame}^- &= \text{pol} \cdot \left(1 - \frac{\beta}{t_{\text{bin}}^+ - t_{\text{bin}}^-}\right) \cdot \mathcal{N} \\ &\quad - \left(\frac{\alpha}{t_{\text{bin}}^+ - t_{\text{bin}}^-}\right) \cdot \text{pol} \cdot \frac{(\mathcal{N}+1)(\mathcal{N})}{2} \\ \text{bin}_{frame}^+ &= \text{pol} \cdot \left(\frac{\beta}{t_{\text{bin}}^+ - t_{\text{bin}}^-}\right) \cdot \mathcal{N} \\ &\quad + \left(\frac{\alpha}{t_{\text{bin}}^+ - t_{\text{bin}}^-}\right) \cdot \text{pol} \cdot \frac{(\mathcal{N}+1)(\mathcal{N})}{2} \end{aligned} \quad (10)$$

These equations are for a pair of frames and are applied for all the pairs in the video sequence to get the full voxel grid. With these closed-form equations and the straight-through-estimator [2] for quantized operations, we can now backpropagate from the binned events inputted to our vision models directly through to the RGC kernels, W , that we want to learn. We additionally extend our formulation to include non-zero refractory periods. See the supplement for these details.

3.2.3. Optimizing Bandwidth and Performance

In this work, we seek to learn the RGC kernels for two tasks: video interpolation and optical flow. We use task-specific loss functions, specifically Charbonnier pixel-wise loss [12, 30] and a masked L1 loss respectively. To optimize for *sparsity* while fitting our learned RGC kernels, we add an additional weighted L1 loss on the number of events to push the model to learn RGC events that maximize performance while minimizing the number of total events.

4. Experiments

We experiment with two tasks to demonstrate the potential benefits of learning RGC events in both DAVIS (intensity + events) and DVS (events only) settings.

1. Video interpolation is one of the most challenging tasks for event cameras and can act as an plug-and-play connection to perception tasks. The state-of-the-art uses output from a DAVIS camera.
2. Optical flow is a particularly useful perception task. In this case, only events are used.

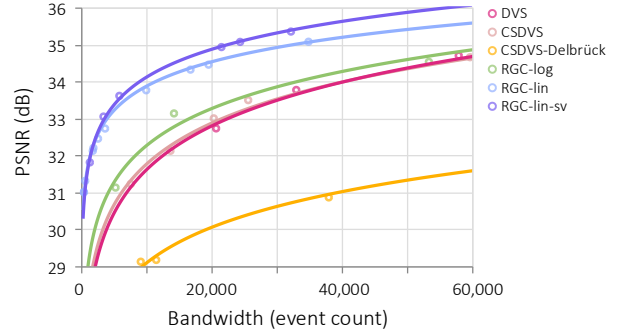


Figure 2. **Video Interpolation Performance vs Bandwidth Trade-off.** We perform video interpolation using DVS, CSDVS, CSDVS-Delbrück, RGC-log (learned, log regime), RGC-lin (learned, linear space), and RGC-lin-sv (learned, linear space, and spatially varying). For any given bandwidth, RGC-lin-sv provides the best performance.

For both tasks, we learn the RGC for improving performance and bandwidth. At the end of this section, we also delve into a number of questions the reader may be curious about, including if learning not one, but multiple kinds of RGC types, improves performance.

4.1. Video Interpolation

We train the state-of-the-art model [56] by Sun et. al that performs video interpolation from events. In the original work, two intensity frames along with the events triggered in-between are fed into their model. From this, they reconstruct 7 frames in-between the two base frames. The original work used ESIM to simulate events from the high-speed GoPro dataset [40]. To learn our events, we replace ESIM with our differentiable simulator. We train with a variety of learned event settings and sparsity weightings to tune the overall average bandwidth. We train the interpolation model with a learning rate of $2e-4$ and our RGC kernels with a learning rate of $5e-5$ end-to-end for 200,000 iterations. We use the Charbonnier pixel-wise loss and our sparsity loss with varying weights to tune performance and bandwidth. Our simulator allows us to learn spatially-varying events too. In this setting, we learn a 2 by 2 bayer-like pattern of kernels and thresholds.

4.2. Optical Flow

For optical flow, we built off of the state-of-the-art optical flow from events model [63] by Wu et. al. We utilize the TartanAir dataset [62] as it has ground truth optical flow. Similarly to [64], we use EMA-VFI [65] to interpolate the RGB video before feeding the frames into our differentiable simulator, interpolating 15 frames between each pair. These intensity frames are used to simulate the RGC events but are not used to recover optical flow. In this task, only events are used to predict flow. We use the ‘‘Hard’’ subset of TartanAir. It provides 18 different scenes, with about 10 tra-

jectories in each scene. Each trajectory has roughly 1,000 frames. To speed up training, we create train, validation, and test datasets from the “Hard” subset of the TartanAir Dataset, saving out random crops of the interpolated images and ground truth optical flow. For maximal generalizability, we split train/val/test at the scene level. The generated dataset has 44,776 training sequences, 5,000 validation sequences, and 5,000 test sequences. Using this dataset and our differentiable event simulator, we learn the RGC kernels and optical flow model end-to-end.

5. Results

5.1. Video Interpolation Results

We present the bandwidth vs. performance tradeoff for interpolating 7 frames between pairs of frames, recovering 240fps from 30fps. Fig. 2 shows the trade-off between bandwidth and performance of different types of events including the traditional DVS, the two variants of the center-surround (CSDVS and the hand-crafted CSDVS_{Delbrück}), our RGC-log (learned RGC in the log intensity domain) and RGC-lin (learned RGC in linear domain). We also show spatially-varying events in the same plot as RGC-lin-sv.

As shown, RGC-lin achieves better performance at any given bandwidth compared to the traditional DVS. Furthermore, adding spatially varying events and thresholds pushes the performance vs bandwidth pareto front up to the left even more. For example, DVS achieves 33.8dB at 33.0k average events per bin while RGC-lin-sv achieves 35.4dB with 32.2k events, a 1.6dB increase using the same number of events. Similarly, if we look at a given performance, such as the 33.8dB that DVS achieves in 33.0k events, RGC-lin achieves 33.8dB with 9.9k events or over 3.3× fewer events, and RGC-lin-sv achieves 33.6dB with 5.9k events or over 5.7× fewer events. These clear benefits highlight the promise of our learned RGC events.

Fig. 3 shows a few examples of the reconstructed videos. The comparisons are for a given average bandwidth of about 20,000 events per bin, specifically 20,716 events for DVS, 20,357 for CSDVS, and 19,557 for RGC-lin. Averaging over the full GoPRO test set sequences, DVS achieves 32.7dB in PSNR, CSDVS achieves 33.0dB, while RGC-lin reaches 34.4dB, 1.67dB higher than DVS. As the middle frame in the sequence is the most challenging to predict, images and metrics (PSNR and SSIM) shown are for the middle frame in each scene. We show the events generated and zoom-ins to details in the reconstructed images. In the ground truth column, we also show the alpha-blended start and end frames of the sequence to illustrate the amount of motion being interpolated. Quantitatively, our learned RGC-lin achieves higher PSNR and SSIM given the same bandwidth as DVS or CSDVS. Qualitatively, the reconstructed structures are sharper and truer to the ground truth.

Table 1. **Optical Flow Quantitative Results.** We compare models trained on DVS with different contrast threshold magnitudes, (which effectively changes the bandwidth), against two of our learned RGC-lin models at different bandwidths. RGC-lin_{lite} and RGC-lin are the same base model, just trained with different sparsity loss weightings. End-point-error (EPE) is the L2 norm between the predicted and ground truth flow. 1PE is the percentage of pixels that have a predicted flow that is off by more than 1 pixel. 3PE is the percentage of pixels off by more than 3 pixels.

	DVS 0.1T	DVS 0.3T	DVS 0.5T	RGC-lin _{lite}	RGC-lin
Bandwidth	7.30M	4.11M	2.77M	2.10M	3.80M
↓ EPE	2.80	3.02	3.33	2.75	2.42
↓ 1PE	55.1	60.8	66.3	52.4	47.1
↓ 3PE	20.1	23.0	26.1	20.7	17.4

5.2. Optical Flow Results

In this task, solely events are used to reconstruct the optical flow. Similarly to interpolation, learning the kernels provides improved performance and lower bandwidth. In Table 1, we compare the models trained end-to-end with our learned RGC-lin to multiple models trained on DVS outputs of different contrast thresholds. The higher the contrast threshold, the fewer the events. RGC-lin and RGC-lin_{lite} are the same model, just trained with different sparsity weightings resulting in different average bandwidths. Over the test set, our learned kernel for optical flow provides the best performance over all metrics. We use the standard metrics: End-point-error (EPE) is the L2-norm between predicted and ground truth flows, 1PE is the percentage of pixels whose flow is off by more than 1 pixel, and 3PE is the percentage of pixels whose flow is off by more than 3 pixels. For DVS, the performance increases with the number of events. However, the best performance comes from the learned RGC-lin kernel, which achieves an EPE of 2.42, better than even the DVS model with nearly twice the number of events. In fact, at the same performance as the DVS 0.1T model, the learned approach only needs 2.10M events, which is 3.5 times fewer events.

Fig. 4 shows three samples of the optical flow reconstruction from DVS and from our Learned RGC-lin events. The DVS model corresponds to DVS 0.1T from tab. 1, as it had the best performance among the DVS models. Learning the kernel end-to-end allows RGC-lin to better reconstruct the flow of the whole frame than DVS while lowering overall bandwidth.

5.3. Insights from the Learned Kernels

In fig. 5, we show a comparison of the DVS, CSDVS, and the learned kernels. As we can see, the best kernels for interpolation and optical flow differ greatly, which is not surprising as the tasks are quite different. However, this highlights how choosing a single sensing kernel like the DVS for every task can limit the potential performance.

For interpolation, the learned kernel reveals that the



Figure 3. **Video Interpolation Qualitative Results.** For each scene, we compare the reconstructions of the middle frame in the sequence for DVS, CSDVS, and RGC-lin. The top row shows the generated events, binned into the corresponding middle time bin, the second is the predicted image and the bottom row shows zoom-ins. The right-most column shows the start and end frames, alpha-blended, ground truth frame, and zoom-ins. PSNR(↑) and SSIM(↑) metrics are shown for each reconstruction.

model places importance on local contrast as its ring structure emphasizes sharp changes in intensity. It effectively captures something akin to an edge detection or bandpass filter, which can be a powerful cue for figuring out how intermediate frames should look.

For optical flow, the kernel has even greater contrast. While the kernel for interpolation was nearly symmetric, the kernel for optical flow is strongly asymmetric, exhibiting some direction-selectivity or the importance of directional gradients. In our eyes, we also have direction-selective RGCs that are more attune to detecting motions, and they are also asymmetric or have elongated receptive fields.

Additionally, in our supplement, we study how the learned kernels change as we increase the receptive field and as we increase bandwidth for video interpolation. In the first study, we sweep kernel sizes $k = 3, 5, 7, 9, 11$ and

find that the performance stays roughly the same at a given bandwidth. In the second study, we increase the weighting on the sparsity loss to tune the bandwidth. We include visualizations of the corresponding kernels in the supplement.

5.4. Multi-Event Exploration

Inspired by the diversity of RGCs in the human eye, we explore models with increasing number of learned RGC types per pixel location for the task of interpolation. With one RGC type, the kernel is roughly symmetric, but as the number of kernels increases, we begin to see asymmetry and direction sensitive features appear. See the supplement for the learned kernels. At a bandwidth of roughly 150,000 events per bin, learning one type of event, RGC_1 , achieves an average PSNR of 36.52dB. Learning just one more kernel, RGC_2 , increases the performance to 36.92dB. With four

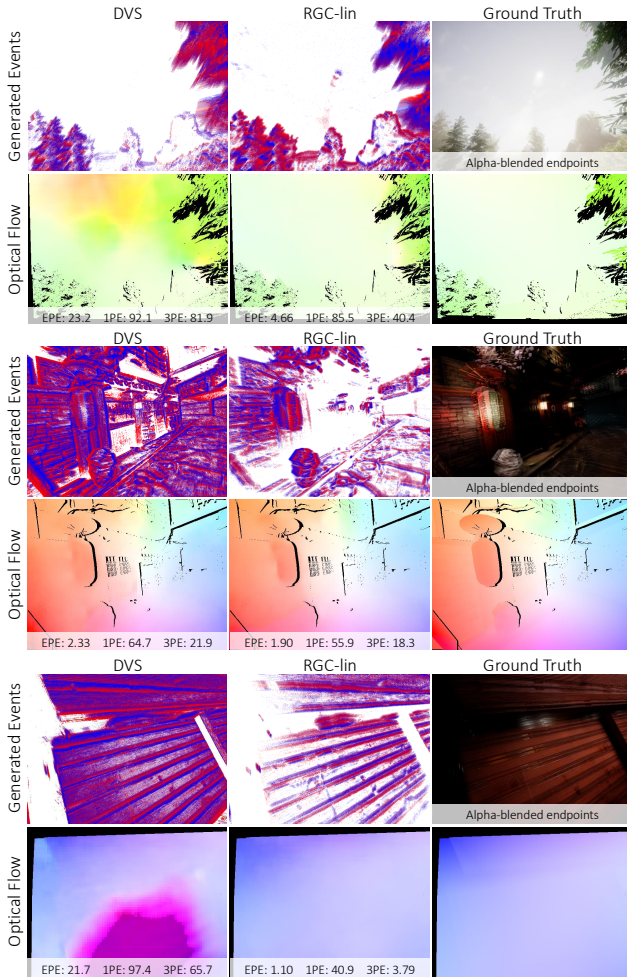


Figure 4. **Optical Flow Qualitative Results.** For each sample, the top row shows the events generated by the DVS kernel and our learned RGC-lin kernel as well as the alpha-blended camera frames *just for reference*. In this task, solely events are used to reconstruct the flow. The bottom row shows the reconstructed flows and the ground truth flow. We show $EPE\downarrow$, $IPE\downarrow$, and $3PE\downarrow$ metrics for the reconstructions.

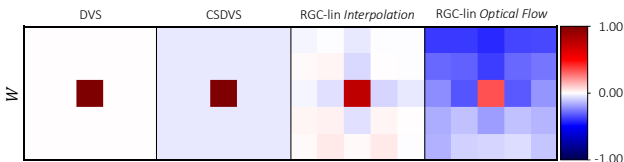


Figure 5. **Comparison of Kernels.** We show the 5×5 kernels for DVS, CSDVS, the learned RGC-lin kernels for video interpolation, and the learned RGC-lin kernels for optical flow.

types of RGCs, RGC_4 , we achieve 37.27dB at the same total bandwidth. In the supplement, we also show a comparison between RGC_1 and RGC_{16} at a lower total bandwidth of 10,000 events. RGC_1 achieves 33.76dB, while RGC_{16} achieves 35.16dB.

6. Discussion and Conclusion

We present a biologically inspired way of sensing by creating a new event-generation framework that extends the traditional pixel-wise event paradigm by leveraging local spatial information. For this purpose, we craft a differentiable event simulator to enable learning of the RGC kernels. We demonstrate experiments on video interpolation and optical flow, showing that our learned events outperform both conventional DVS and center-surround DVS (CSDVS) methods under the same bandwidth constraints. By utilizing neighborhood context, our proposed approach delivers richer, more informative event streams—pointing to a promising new direction in event-based sensing for real-time and resource-constrained applications.

Hardware Feasibility. There are several avenues for creating Neural Ganglion Sensors in hardware, especially with the rise of in-pixel compute. [17] explored the feasibility of a center-surround DVS and proposed using polysilicon lateral resistors to weight the surround according to their hand-crafted kernel. A similar resistor mesh could be applied with our learned kernel weights. For more adaptability where the weights can be dynamically updated on-the-fly, emerging in-pixel compute platforms such as [11] can also be an attractive candidate. These emerging platforms integrate compute and memory in the sensor plane and have already shown promise in a number of computational imaging pipelines [37, 42, 53]. Implementing the ideal learned RGC kernels requires a multiplication in floating point precision. However, the operations available on in-pixel processors are currently limited. Some approximations or training with additional constraints would be required for implementation on the current generation of sensor-processors, such as operating in highly quantized regimes where kernels are binary or ternary [4, 5, 34, 54]. Furthermore, [6] has demonstrated multiple highly quantized filters per pixel, illustrating a potential route to multi-channel RGC events. As the capabilities of in-pixel compute continue to develop, our full floating point RGC kernels will become feasible.

Future Directions. There are many potential extensions to this work. In particular, exploring non-binary nonlinearities could be beneficial to the vision tasks. Thresholding is currently used as it best aligns with what we see in existing event sensors, and the binary spikes are akin to the spiking potential in human retinas.

As the demand for efficient sensing continues to grow across domains, including for AR/VR headsets and drones, we hope these insights will help guide sensor designers in optimizing future event-sensing technologies. With more compute and memory becoming available in emerging sensor-processors, Neural Ganglion Sensors will provide a promising balance: a minimal amount of on-sensor compute for massive bandwidth savings.

Acknowledgments

This project was in part supported by Samsung, SK Hynix, and the NSF Graduate Research Fellowship Program.

References

- [1] Arnon Amir, Brian Taba, David Berg, Timothy Melano, Jeffrey McKinstry, Carmelo Di Nolfo, Tapan Nayak, Alexander Andreopoulos, Guillaume Garreau, Marcela Mendoza, Jeff Kunitz, Michael Debole, Steve Esser, Tobi Delbruck, Myron Flickner, and Dharmendra Modha. A low power, fully event-based gesture recognition system. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7388–7397, 2017. 2
- [2] Yoshua Bengio, Nicholas Léonard, and Aaron C. Courville. Estimating or propagating gradients through stochastic neurons for conditional computation. *ArXiv*, abs/1308.3432, 2013. 5
- [3] Raphael Berner, Christian Brandli, Minhao Yang, Shih-Chii Liu, and Tobi Delbruck. A 240x180 120db 10mw 12us-latency sparse output vision sensor for mobile applications. 2013. 2
- [4] L. Bose, P. Dudek, J. Chen, S. Carey, and W. Mayol-Cuevas. A camera that cnns: Towards embedded neural networks on pixel processor arrays. In *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 1335–1344, Los Alamitos, CA, USA, 2019. IEEE Computer Society. 8
- [5] Laurie Bose, Piotr Dudek, Jianing Chen, Stephen J. Carey, and Walterio W. Mayol-Cuevas. Fully embedding fast convolutional networks on pixel processor arrays. In *Computer Vision – ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXIX*, page 488–503, Berlin, Heidelberg, 2020. Springer-Verlag. 8
- [6] Laurie Bose, Piotr Dudek, Stephen J. Carey, and Jianing Chen. Live demonstration: Scamp-7. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, pages 3995–3996, 2023. 8
- [7] Vicente Botella-Soler, Stéphane Deny, Georg Martius, Olivier Marre, and Gašper Tkačik. Nonlinear decoding of a complex movie from the mammalian retina. *PLOS Computational Biology*, 14(5):1–27, 2018. 2
- [8] János Botzheim, Takenori Obo, and Naoyuki Kubota. Human gesture recognition for robot partners by spiking neural network and classification learning. In *The 6th International Conference on Soft Computing and Intelligent Systems, and The 13th International Symposium on Advanced Intelligence Systems*, pages 1954–1958, 2012. 2
- [9] Nora Brackbill, Colleen Rhoades, Alexandra Kling, Nishal P Shah, Alexander Sher, Alan M Litke, and EJ Chichilnisky. Reconstruction of natural images from responses of primate retinal ganglion cells. *eLife*, 9:e58516, 2020. 2
- [10] Christian Brandli, Raphael Berner, Minhao Yang, Shih-Chii Liu, and Tobi Delbruck. A 240 × 180 130 db 3 μs latency global shutter spatiotemporal vision sensor. *IEEE Journal of Solid-State Circuits*, 49(10):2333–2341, 2014. 2
- [11] Stephen J. Carey, Alexey Lopich, David R.W. Barr, Bin Wang, and Piotr Dudek. A 100,000 fps vision sensor with embedded 535gops/w 256×256 simd processor array. In *2013 Symposium on VLSI Circuits*, pages C182–C183, 2013. 8
- [12] P. Charbonnier, L. Blanc-Feraud, G. Aubert, and M. Barlaud. Two deterministic half-quadratic regularization algorithms for computed imaging. In *Proceedings of 1st International Conference on Image Processing*, pages 168–172 vol.2, 1994. 5
- [13] EJ Chichilnisky. A simple white noise analysis of neuronal light responses. In *Network (Bristol, England)*, pages 199–213, 2001. 2, 3
- [14] E. Culurciello, R. Etienne-Cummings, and K.A. Boahen. A biomorphic digital image sensor. *IEEE Journal of Solid-State Circuits*, 38(2):281–294, 2003. 2
- [15] Thomas Dalgaty, Thomas Mesquida, Damien Joubert, Amos Sironi, Cyrille Soubeyrat, Pascal Vivet, and Christoph Posch. The cnn vs. snn event-camera dichotomy and perspectives for event-graph neural networks. In *2023 Design, Automation & Test in Europe Conference & Exhibition (DATE)*, pages 1–6, 2023. 2
- [16] Tobi Delbruck. Neuromorphic vision sensing and processing. *2016 46th European Solid-State Device Research Conference (ESSDERC)*, pages 7–14, 2016. 2
- [17] Tobi Delbruck, Chenghan Li, Rui Graca, and Brian Mcreynolds. Utility and feasibility of a center surround event camera. 2022. 2, 4, 8
- [18] Yongjian Deng, Youfu Li, and Hao Chen. Amae: Adaptive motion-agnostic encoder for event-based object classification. *IEEE Robotics and Automation Letters*, 5(3):4596–4603, 2020. 2
- [19] G.D. Field and E.J. Chichilnisky. Information processing in the primate retina: Circuitry and coding. *Annual Review of Neuroscience*, 30(1):1–30, 2007. PMID: 17335403. 2
- [20] G. Gallego, T. Delbruck, G. Orchard, C. Bartolozzi, B. Taba, A. Censi, S. Leutenegger, A. J. Davison, J. Conradt, K. Daniilidis, and D. Scaramuzza. Event-based vision: A survey. *IEEE Transactions on Pattern Analysis & Machine Intelligence*, 44(01):154–180, 2022. 1, 2
- [21] Yue Gao, Jiaxuan Lu, Siqi Li, Nan Ma, Shaoyi Du, Yipeng Li, and Qionghai Dai. Action recognition and benchmark using event cameras. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(12):14081–14097, 2023. 1
- [22] Rohan Ghosh, Anupam Kumar Gupta, Andrei Nakagawa Silva, Alcimar Barbosa Soares, and Nitish V. Thakor. Spatiotemporal filtering for event-based action recognition. *ArXiv*, abs/1903.07067, 2019. 1
- [23] Fuqiang Gu, Weicong Sng, Tasbolat Taunyazov, and Harold Soh. Tactilesnet: A spiking graph neural network for event-based tactile object recognition. In *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2020. 2
- [24] Jinjin Gu, Jinan Zhou, Ringo Sai Wo Chu, Yan Chen, Jiawei Zhang, Xuanye Cheng, Song Zhang, and Jimmy S. Ren. Self-supervised intensity-event stereo matching. 2022. 1
- [25] Tianruo Guo, David Tsai, Siwei Bai, John Morley, Gregg Suaning, Nigel Lovell, and Socrates Dokos. Understanding

- the retina: A review of computational models of the retina from the single cell to the network level. *Critical reviews in biomedical engineering*, 42:419–36, 2014. 2
- [26] A.L. Hodgkin and A.F. Huxley. A quantitative description of membrane current and its application to conduction and excitation in nerve. *Bulletin of Mathematical Biology*, 52(1): 25–71, 1990. 2
- [27] Y Hu, S C Liu, and T Delbruck. v2e: From video frames to realistic DVS events. In *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*. IEEE, 2021. 4
- [28] Renaud Jolivet, Timothy J., and Wulfram Gerstner. The spike response model: A framework to predict neuronal spike trains. In *Artificial Neural Networks and Neural Information Processing — ICANN/ICONIP 2003*, pages 846–853, Berlin, Heidelberg, 2003. Springer Berlin Heidelberg. 2
- [29] Stephen W Kuffler. Discharge patterns and functional organization of mammalian retina. *Journal of neurophysiology*, 16(1):37–68, 1953. 1
- [30] Wei-Sheng Lai, Jia-Bin Huang, Narendra Ahuja, and Ming-Hsuan Yang. Deep laplacian pyramid networks for fast and accurate super-resolution. In *IEEE Conferene on Computer Vision and Pattern Recognition*, 2017. 5
- [31] Chenghan Li. Two-stream vision sensors. 2017. 2
- [32] Patrick Lichtsteiner, Christoph Posch, and Tobi Delbruck. A 128×128 120 db 15 μ s latency asynchronous temporal contrast vision sensor. *IEEE Journal of Solid-State Circuits*, 43(2):566–576, 2008. 2
- [33] Siying Liu and Pier Luigi Dragotti. Sensing diversity and sparsity models for event generation and video reconstruction from events. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(10):12444–12458, 2023. 2
- [34] Yanan Liu, Laurie Bose, Jianing Chen, Stephen J. Carey, Piotr Dudek, and W. Mayol-Cuevas. High-speed light-weight cnn inference via strided convolutions on a pixel processor array. In *British Machine Vision Conference*, 2020. 8
- [35] Misha Mahowald. *VLSI analogs of neuronal visual processing: a synthesis of form and function*. 2008. 2
- [36] Ana María Iglesias Maqueda, Antonio Loquercio, Guillermo Gallego, Narciso García, and Davide Scaramuzza. Event-based vision meets deep learning on steering prediction for self-driving cars. *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5419–5427, 2018. 2
- [37] Julien N. P. Martel, Lorenz K. Müller, Stephen J. Carey, Piotr Dudek, and Gordon Wetzstein. Neural sensors: Learning pixel exposures for hdr imaging and video compressive sensing with programmable sensors. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 42(7):1642–1653, 2020. 2, 8
- [38] Lane T. McIntosh, Niru Maheswaranathan, Aran Nayebi, Surya Ganguli, and Stephen A. Baccus. Deep learning models of the retinal response to natural scenes. In *Proceedings of the 30th International Conference on Neural Information Processing Systems*, page 1369–1377, Red Hook, NY, USA, 2016. Curran Associates Inc. 2
- [39] Manasi Muglikar, Leonard Bauersfeld, Diederik Moeys, and Davide Scaramuzza. Event-based shape from polarization. In *IEEE / CVF Computer Vision and Pattern Recognition Conference (CVPR)*, 2023. 1
- [40] Seungjun Nah, Tae Hyun Kim, and Kyoung Mu Lee. Deep multi-scale convolutional neural network for dynamic scene deblurring. In *CVPR*, 2017. 5
- [41] Shree K. Nayar and Tomoo Mitsunaga. High dynamic range imaging: spatially varying pixel exposures. *Proceedings IEEE Conference on Computer Vision and Pattern Recognition. CVPR 2000 (Cat. No.PR00662)*, 1:472–479 vol.1, 2000. 2
- [42] Cindy M Nguyen, Julien NP Martel, and Gordon Wetzstein. Learning spatially varying pixel exposures for motion deblurring. *IEEE International Conference on Computational Photography (ICCP)*, 2022. 2, 8
- [43] Nikhil Parthasarathy, Eleanor Batty, William Falcon, Thomas Rutten, Mohit Rajpal, E.J. Chichilnisky, and Liam Paninski. Neural networks for efficient bayesian decoding of natural images from retinal neurons. In *Advances in Neural Information Processing Systems*. Curran Associates, Inc., 2017. 2
- [44] Jonathan W. Pillow, Liam Paninski, Valerie J. Uzzell, Eero P. Simoncelli, and E. J. Chichilnisky. Prediction and decoding of retinal ganglion cell responses with a probabilistic spiking model. *Journal of Neuroscience*, 25(47):11003–11013, 2005. 2
- [45] Chiara Plizzari, Mirco Planamente, Gabriele Goletto, Marco Cannici, Emanuele Gusso, Matteo Matteucci, and Barbara Caputo. E2(go)motion: Motion augmented event stream for egocentric action recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 19935–19947, 2022. 1
- [46] Christoph Posch, Daniel Matolin, and Rainer Wohlgenannt. A qvga 143 db dynamic range frame-free pwm image sensor with lossless pixel-level video compression and time-domain cds. *IEEE Journal of Solid-State Circuits*, 46(1):259–275, 2011. 2
- [47] Christoph Posch, Teresa Serrano-Gotarredona, Bernabe Linares-Barranco, and Tobi Delbruck. Retinomorphc event-based vision sensors: Bioinspired cameras with spiking output. *Proceedings of the IEEE*, 102(10):1470–1484, 2014. 2
- [48] Henri Rebecq, René Ranftl, Vladlen Koltun, and Davide Scaramuzza. High speed and high dynamic range video with an event camera. *IEEE Trans. Pattern Anal. Mach. Intell. (T-PAMI)*, 2019. 2
- [49] B Roska and M Meister. The retina dissects the visual scene into distinct features. In *The New Visual Neurosciences (Werner, JS, Chalupa, LM, eds)*, pages 163–182, 2014. 2, 4
- [50] Joshua Sanes and Richard Masland. The types of retinal ganglion cells: Current status and implications for neuronal classification. *Annual review of neuroscience*, 38, 2015. 1
- [51] Shintaro Shiba, Yoshimitsu Aoki, and Guillermo Gallego. Secrets of event-based optical flow. In *European Conference on Computer Vision (ECCV)*, pages 628–645, 2022. 1
- [52] Shintaro Shiba, , Friedhelm Hamann, Yoshimitsu Aoki, and Guillermo Gallego. Event-based background-oriented schlieren. 2023. 1

- [53] H. M. So, J. P. Martel, G. Wetzstein, and P. Dudek. Mantisacam: Learning snapshot high-dynamic-range imaging with perceptually-based in-pixel irradiance encoding. In *2022 IEEE International Conference on Computational Photography (ICCP)*, pages 1–12, Los Alamitos, CA, USA, 2022. IEEE Computer Society. 8
- [54] Haley M. So, Laurie Bose, Piotr Dudek, and Gordon Wetzstein. Pixelrnn: In-pixel recurrent neural networks for end-to-end-optimized perception with neural sensors. In *CVPR*, pages 25233–25244, 2024. 8
- [55] Qiaoyi Su, Yuhong Chou, Yifan Hu, Jianing Li, Shijie Mei, Ziyang Zhang, and Guoqi Li. Deep directly-trained spiking neural networks for object detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6555–6565, 2023. 2
- [56] Lei Sun, Christos Sakaridis, Jingyun Liang, Peng Sun, Jiezhong Cao, Kai Zhang, Qi Jiang, Kaiwei Wang, and Luc Van Gool. Event-based frame interpolation with ad-hoc deblurring. *arXiv preprint arXiv:2301.05191*, 2023. 4, 5
- [57] Varun Sundar, Matt Dutton, Andrei Ardelean, Claudio Bruschi, Edoardo Charbon, and Mohit Gupta. Generalized event cameras. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024. 2
- [58] Yi Tian and Juan Andrade-Cetto. Event transformer flownet for optical flow estimation. In *33rd British Machine Vision Conference 2022, BMVC 2022, London, UK, November 21-24, 2022*. BMVA Press, 2022. 1
- [59] Todd W. Troyer and Kenneth D. Miller. Physiological Gain Leads to High ISI Variability in a Simple Model of a Cortical Regular Spiking Cell. *Neural Computation*, 9(5):971–983, 1997. 2
- [60] Stepan Tulyakov, Francois Fleuret, Martin Kiefel, Peter Gehler, and Michael Hirsch. Learning an event sequence embedding for dense event-based deep stereo. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019. 1
- [61] Brian A Wandell. *Foundations of vision*. Sinauer Associates, 1995. 1
- [62] Wenshan Wang, Delong Zhu, Xiangwei Wang, Yaoyu Hu, Yuheng Qiu, Chen Wang, Yafei Hu, Ashish Kapoor, and Sebastian Scherer. Tartanair: A dataset to push the limits of visual slam. 2020. 5
- [63] Yilun Wu, Federico Paredes-Vallés, and Guido C. H. E. de Croon. Lightweight event-based optical flow estimation via iterative deblurring. In *Proceedings of IEEE International Conference on Robotics and Automation (ICRA'24)*, 2024. To Appear. 4, 5
- [64] Yan Yang, Liyuan Pan, and Liu Liu. Event camera data pre-training. *CoRR*, abs/2301.01928, 2023. 5
- [65] Guozhen Zhang, Yuhan Zhu, Haonan Wang, Youxin Chen, Gangshan Wu, and Limin Wang. Extracting motion and appearance via inter-frame attention for efficient video frame interpolation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5682–5692, 2023. 5
- [66] Xu Zheng, Yexin Liu, Yunfan Lu, Tongyan Hua, Tianbo Pan, Weiming Zhang, Dacheng Tao, and Lin Wang. Deep learning for event-based vision: A comprehensive survey and benchmarks. *arXiv preprint arXiv:2302.08890*, 2023. 2
- [67] Xueye Zheng, Ye-Peng Liu, Yunfan Lu, Tongyan Hua, Tianbo Pan, Weiming Zhang, Dacheng Tao, and Lin Wang. Deep learning for event-based vision: A comprehensive survey and benchmarks. *ArXiv*, abs/2302.08890, 2023. 2
- [68] Zhaokun Zhou, Yuesheng Zhu, Chao He, Yaowei Wang, Shuicheng YAN, Yonghong Tian, and Li Yuan. Spikformer: When spiking neural network meets transformer. In *The Eleventh International Conference on Learning Representations*, 2023. 2
- [69] Alex Zhu, Liangzhe Yuan, Kenneth Chaney, and Kostas Daniilidis. Ev-flownet: Self-supervised optical flow estimation for event-based cameras. In *Proceedings of Robotics: Science and Systems*, Pittsburgh, Pennsylvania, 2018. 2
- [70] Lin Zhu, Xiao Wang, Yi Chang, Jianing Li, Tiejun Huang, and Yonghong Tian. Event-based video reconstruction via potential-assisted spiking neural network. In *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3584–3594, 2022. 2