

# EVENTSR-ZERO: TRAINING-FREE EVENT VIDEO SUPER-RESOLUTION WITH DIFFUSION PRIORS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Event-to-Video (E2V) methods aim to reconstruct intensity frames from events, bridging the gap between event-based and image-based vision. However, existing E2V approaches often fail to recover fine structures, leading to reconstructions with artifacts and degraded quality. To address this, we explore the task of Event-to-Video Super-Resolution (EVSR), which aims to reconstruct high-resolution video from low-resolution events. We present EventSR-Zero, a training-free framework that exploits the high temporal resolution of event cameras to recover fine-grained details from low-resolution events and uses them to guide a diffusion-based Video Super-Resolution (VSR) model in generating high-quality super-resolved videos of the underlying scene. Our approach incorporates two key components: (1) an Implicit Contrast Refinement (ICR) module that robustly extracts sub-pixel scene details from low-resolution events, and (2) a Reconditioning Guidance (RG) module that reliably steers the diffusion VSR process using the high-resolution event signal from ICR. Extensive experiments demonstrate that EventSR-Zero achieves state-of-the-art performance, surpassing existing event-based super-resolution methods. We will release our source code upon acceptance.

## 1 INTRODUCTION

Event cameras represent a significant advancement in vision technology, addressing the limitations of traditional frame-based cameras in high-speed and high-dynamic-range scenarios. Unlike conventional cameras that capture entire frames at fixed intervals, event cameras operate asynchronously, detecting per-pixel brightness changes with microsecond precision. Their sub-millisecond resolution, wide dynamic range, and sparse, memory-efficient output enable high-speed, low-latency vision while reducing storage and computational demands. Event cameras are crucial in robotics and autonomous systems for real-time navigation and obstacle detection Falanga et al. (2020); Forrai et al. (2023); Huang et al. (2024). They also enhance object recognition Zubić et al. (2024); Gehrig & Scaramuzza (2023); Li et al. (2021); Mitrokhin et al. (2019); Scheerlinck et al. (2018) and surveillance Bi et al. (2019); Freeman et al. (2024); Verma et al. (2024), particularly in low-light conditions. In AR/VR, they enable precise gesture tracking and interaction Plizzari et al. (2022); Gao et al. (2023; 2024).

Event-to-video (E2V) reconstruction is crucial for bridging the gap between event-based and conventional vision systems, enabling broader adoption of event cameras in existing applications. Since event cameras output sparse, asynchronous data that encode only brightness changes, they lack absolute intensity information which makes direct interpretation challenging. Reconstructing videos from events restores intensity frames, allowing seamless integration with traditional computer vi-

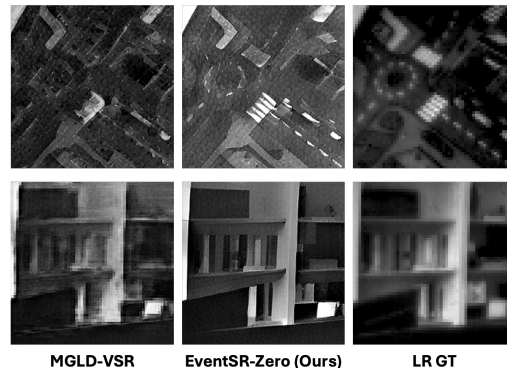


Figure 1: Our EventSR-Zero guides a diffusion VSR model with high-resolution structures recovered from events to generate high-quality SR frames from low-resolution events.

054 sion pipelines that rely on frame-based processing. E2V reconstruction also enhances applications  
055 such as high-speed video capture, where traditional cameras struggle with motion blur, and low-light  
056 imaging, where event cameras excel due to their high dynamic range.

057 Research on deep-learning based E2V reconstruction Rebecq et al. (2019); Scheerlinck et al. (2020);  
058 Stoffregen et al. (2020a); Paredes-Valles & de Croon (2021); Cadena et al. (2021); Weng et al.  
059 (2021); Mostafavi et al. (2020); Wang et al. (2020); Ercan et al. (2024) aims to convert raw event  
060 streams into video frames by training over paired event-image datasets. This transformation enables  
061 more intuitive and accessible representations of the rich temporal and structural information con-  
062 tained within sparse, unstructured event data. However, despite significant progress, existing E2V  
063 techniques still struggle to recover fine details, particularly in complex scenes or high-speed scenar-  
064 ios, where the inherent sparsity of event data and the absence of absolute intensity values lead to  
065 artifact-laden or low-quality reconstructions.

066 Building on the limitations of conventional Event-to-Video (E2V) methods, Event-to-Video Super-  
067 Resolution (EVSR) introduces a significant advancement by enhancing the spatial resolution of  
068 reconstructed frames. Unlike traditional cameras that capture purely spatial information, event cam-  
069 eras record a spatiotemporal stream with microsecond temporal precision. This high temporal res-  
070 olution enables the capture of subtle displacements that provide complementary spatial cues over  
071 time, thereby allowing inference of details beyond the native resolution of the sensor. Leveraging  
072 this property, EVSR can reconstruct high-frequency scene structures and generate super-resolved  
073 intensity images directly from event data.

074 By incorporating sub-pixel event alignment strategies with deep image priors, EVSR can signifi-  
075 cantly improve image sharpness and edge definition, making event-based vision more practical for  
076 high-precision applications such as autonomous navigation, medical imaging, and high-speed video  
077 capture, where fine visual details are crucial. Moreover, achieving high-resolution reconstructions  
078 from events enables better compatibility with existing high-resolution vision models, expanding  
079 the applicability of event cameras in mainstream computer vision tasks. The development of ro-  
080 bust EVSR techniques can unlock the full potential of event cameras, enabling not just fast and  
081 low-latency vision but also high-quality, detailed imagery suitable for a wide range of real-world  
082 applications.

083 A direct approach is to train a high-resolution (HR) event-to-video model to learn the SR prior in  
084 event space. However, real-world HR event datasets do not exist, and furthermore existing strategies  
085 to simulate HR event datasets are still encumbered by a sim-to-real domain gap. As a result, EVSR  
086 remains a complex, ongoing challenge in event-based vision research.

087 In this work, we propose EventSR-Zero, a training-free approach that robustly recovers fine-grained  
088 details from low-resolution events by exploiting the high temporal resolution of event cameras and  
089 uses these details to guide a diffusion-based VSR model to generate highly detailed frames from  
090 low-resolution event streams.

091 Our method comprises two key components:

092 The **Implicit Contrast Refinement (ICR)** module enables recovery of high-frequency sub-pixel  
093 details from LR event data. It formulates a high-resolution contrast maximization (CMax) space,  
094 regularized by a frequency-constrained lightweight MLP that serves as an implicit function to miti-  
095 gate the high propensity of event collapse in HR space.

096  
097 2) The **Reconditioning Guidance (RG)** module is a novel diffusion guidance strategy that steers the  
098 diffusion trajectory through controlled adjustments to the conditioning image. At each step, RG de-  
099 rives frame estimates from intermediate latents and aligns their flow-directed spatial gradients with  
100 the HR details recovered by ICR, thereby transferring fine event structures into the conditioning.  
101 This ensures that the diffusion process generates frames whose structural features remain consistent  
102 with the high-resolution event details provided by ICR.

103 Our method does not utilize task-specific training data, thus alleviating the requirement for high-  
104 resolution event-image datasets. Qualitative and quantitative results show that our EventSR-Zero  
105 outperforms existing baselines to achieve effective event video SR. Our **main contributions** are as  
106 follows:

107

- We introduce EventSR-Zero, a training-free method that achieves high quality super-resolved intensity images from LR events.
- We propose Implicit Contrast Refinement (ICR) that mitigates event collapse in HR space to reliably recover sub-pixel level details from high-temporal resolution events.
- We propose Reconditioning Guidance (RG), a novel diffusion guidance technique that produces event-guided, super-resolved intensity frames while maintaining consistency towards underlying HR scene structures.

## 2 RELATED WORKS

### 2.1 EVENT TO VIDEO SUPER-RESOLUTION (EVSR)

Reconstructing intensity images from events is a key topic in event-based vision, with various methods offering different assumptions and processing techniques. Early approaches Cook et al. (2011); Kim et al. (2014); Agrawal et al. (2005); Kim et al. (2016); Barua et al. (2016); Aharon et al. (2006); Bardow et al. (2016); Munda et al. (2018); Scheerlinck et al. (2018) relied on simplified assumptions in constrained camera motion or brightness constancy to reconstruct intensity frames. More recently, deep learning methods Rebecq et al. (2019); Scheerlinck et al. (2020); Stoffregen et al. (2020a); Paredes-Valles & de Croon (2021); Cadena et al. (2021); Weng et al. (2021); Mostafavi et al. (2020); Wang et al. (2020); Ercan et al. (2024) have achieved state of the art results in intensity image reconstruction. These methods adopt voxelized event grids to encode sparse events, and typically operate a recurrent network to capture long-range context from past event segments. These methods are trained over low-resolution event-image datasets and are effective at video reconstruction at low-resolution. Across existing methods, Hyper-E2VID Ercan et al. (2024) currently achieves the highest reconstruction fidelity through the combination of a recurrent event voxel encoding architecture and a dynamic filter generation hypernetwork.

In the area of Event to Video Super-Resolution (EVSR), super-resolved intensity frames are produced from pure event streams. Mostafavi et al. (2020) parses events as stacked event images and implements an optical flow estimator, a feature rectification network, a recurrent SR network, and a mixer network to reconstruct HR frames. The modules are trained end-to-end over a ESIM Rebecq et al. (2018) simulated HR dataset to achieve intensity image SR. Wang et al. (2020) proposes three sequential networks for reconstruction, restoration and SR that are trained end-to-end over a simulated EventSR dataset that is also generated from ESIM Rebecq et al. (2018). Duan et al. (2021) proposes a display-camera system for HR event data collection, which is used to train a U-Net based framework to estimate HR spatiotemporal event point clouds.

Similarly, our work seeks to produce super-resolved event-based intensity images solely from event streams. Unlike existing approaches, our EventSR-Zero does not require synthetic HR event datasets. In our work, we use Hyper-E2VID Ercan et al. (2024) to generate the initial LR frames for event-guided VSR. We also adopt the formulation of Zhang et al. (2023) to provide event-based guidance on the diffusion trajectory by aligning the spatial gradients of the intermediate image estimates with the Image of Warped Events (IWE) that contain HR details captured from ICR.

### 2.2 DIFFUSION-BASED VIDEO SUPER-RESOLUTION

Diffusion models transform a sample from a noise distribution to a target data distribution using a fixed forward process and a learned reverse denoising process. The reverse process is learned by a network  $\phi$  that is trained to denoise a noised latent  $\mathbf{z}_t$  by predicting its noise component  $\epsilon_\phi(\mathbf{z}_t, y, t)$  conditioned upon  $y$  (text/image) and timestep  $t$ . Diffusion-based Video Super-Resolution (VSR) models take a sequence of LR images as conditioning for the diffusion network  $\phi$  and adds SR details to the generated image across multiple diffusion steps. The reverse process adds detail to the generated image by iteratively denoising a noise sample  $\mathbf{z}_T$  into a sample  $\mathbf{z}_0$  from the data distribution across a diffusion trajectory.

Recently, diffusion-based single-image SR models Rombach et al. (2022); Wang et al. (2024) have been adapted for Video Super-Resolution (VSR) Yang et al. (2024); Zhou et al. (2024), incorporating fine-tuning and specific constraints to manage temporal coherence and motion consistency across frames. These adaptations allow diffusion-based VSR models to set new benchmarks in video super-

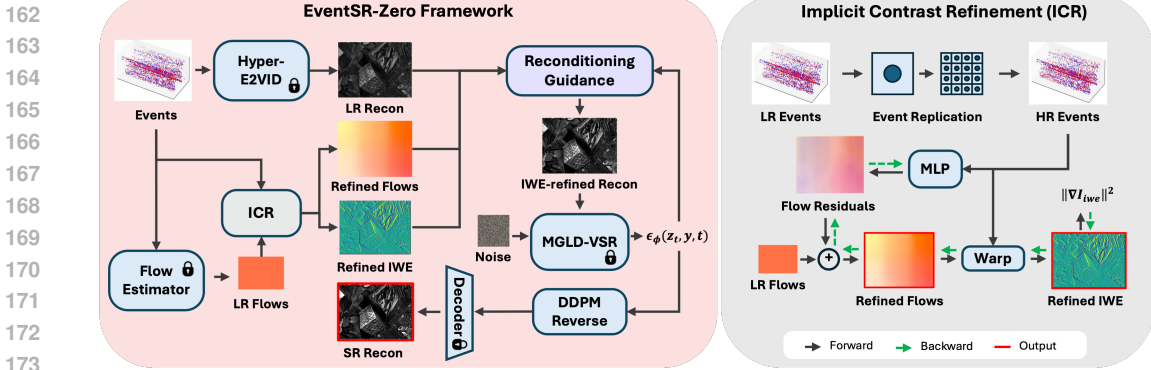


Figure 2: (Left) EventSR-Zero enhances LR reconstructions from HyperE2VID using a diffusion VSR model (MGLD-VSR) guided by Conditioning Guidance (RG), which aligns latent reconstructions with HR IWEs and flows from our Implicit Contrast Refinement (ICR) module. (Right) ICR replicates events over subpixels and encodes HR coordinates with frequency-controlled positional embeddings, passing them through a lightweight MLP to estimate smooth HR flow residuals.

resolution quality. In this work, we utilize a frozen MGLD-VSR Yang et al. (2024) model guided by motion-compensated events to achieve high-quality event-based video super-resolution.

### 2.3 CONTRAST MAXIMIZATION

Contrast maximization (CMax) is a technique in event-based vision that optimizes flow estimates for motion compensation, improving event alignment and enhancing scene edges. It has been applied to motion estimation Gallego et al. (2018) and intensity reconstruction Zhang et al. (2023), where contrast-maximized stacks sharpen visual features. However, CMax is prone to event collapse Shiba et al. (2022a), where events aggregate into patterns that achieve high contrast but obscure true structure. To mitigate this, recent works impose hierarchical grids that interpolate pixel flows from coarser estimates Shiba et al. (2022b), or derive flow via nearest-neighbor associations to nonlinear motion trajectories Friedhelm Hamann (2024). Others reduce collapse risk by restricting motion model degrees of freedom, e.g., assuming rotational or 6-DOF egomotion Zhang et al. (2023). While these strategies improve stability, they oversimplify real-world scene dynamics and cannot capture complex non-rigid motions. In contrast, EventSR-Zero departs from prior work by applying CMax in a high-resolution *refinement* stage rather than at the native event resolution. Specifically, EventSR-Zero introduces a high-resolution optimization space that simulates events from a finer-resolution sensor, allowing convergence to more detailed scene structures. To address the heightened risk of event collapse in this enlarged space, it employs a frequency-constrained implicit function regularizer that stabilizes the optimization.

## 3 OUR METHOD

Our EventSR-Zero framework integrates an event-to-video reconstruction model Hyper-E2VID Ercan et al. (2024) with a diffusion-based video super-resolution (VSR) model MGLD-VSR Yang et al. (2024) to create a backbone for event video super-resolution. Hyper-E2VID initially reconstructs low-resolution (LR) frames from raw event streams, which are then super-resolved into high-resolution (HR) frames by MGLD-VSR. However, the VSR model cannot differentiate between artifacts and true features in the LR input. Without proper guidance, the ill-posed nature of super-resolution where multiple SR solutions can correspond to a single LR input causes this pipeline to risk generating details that deviate from the true scene structure.

To address the above-mentioned issues, we introduce a training-free event-informed guidance that tweaks the diffusion trajectory to help produce details that are aligned with HR structural features derived from the event stream. Our approach includes two main components: 1) An **Implicit Contrast Refinement (ICR)** module to recover high-resolution details from events; 2) A **Conditioning Guidance (RG)** module to provide precise diffusion trajectory adjustments for the VSR model.

Fig. 2 (left) illustrates the EventSR-Zero framework. Raw events are processed by Hyper-E2VID Ercan et al. (2024) and an event-based flow estimator (Shiba et al. (2022b)) to produce initial LR reconstructions and LR flows, respectively. Our ICR then upscales the LR flows to HR flows, producing refined IWEs to assist in diffusion guidance. The MGLD-VSR model Yang et al. (2024) takes the LR reconstructions from Hyper-E2VID as conditioning inputs for the diffusion-based super-resolution process. At each diffusion step, the Reconstruction Guidance (RG) module extracts intermediate reconstructions from the diffusion latents and aligns their spatial gradients with the refined IWEs, which embed high-frequency scene details. These enriched reconstructions are iteratively recombined into the conditioning images, guiding subsequent diffusion steps to generate super-resolved frames with improved fidelity to fine-grained structures captured by the event data.

### 3.1 IMPLICIT CONTRAST REFINEMENT

Fig. 2 (right) outlines Implicit Contrast Refinement (ICR), which begins with event replication: each pixel in the LR event sensor is subdivided to simulate a higher-resolution sensor. For  $4\times$  upscaling, every LR pixel corresponds to 16 subpixels in HR space. This is achieved by replicating each LR event across the equivalent patch of HR pixels, with each HR event given its respective spatial offset. This enables the refinement process to consolidate details in a higher spatial resolution. This event replication crucial for initializing the ICR process is formulated as:

$$e_{x,y,t,p}^{lr} = \{e_{kx,ky,t,p}^{hr}, \dots, e_{kx+(k-1),ky+(k-1),t,p}^{hr}\} \quad (1)$$

where  $x, y, t, p$  is the location, time and polarity, and  $k = 4$  is the scaling factor for  $4\times$  SR. We run Contrast Maximization (CMax) as a refinement step in HR space to recover HR flow solutions that reveal finer scene structures from the LR event stream. This is done by finding HR flow estimates that maximize the alignment of corresponding events. Specifically, a warp function  $\mathbf{f}(\cdot)$  maps each event in a set  $\mathcal{E} = \{e_{x,y,t,p}^i\}_{i=1}^{N_e}$  to a new position:

$$e_{x',y',t',p} = \mathbf{f}(e_{x,y,t,p}, t'). \quad (2)$$

The warped events are accumulated onto an Image of Warped Events (IWE) defined as:

$$I_{iwe}(\mathbf{p}; \mathbf{f}) \doteq \sum_{i=1}^{N_e} p_k \delta(\mathbf{p} - \mathbf{f}(\mathbf{p}'_i)), \quad (3)$$

where each pixel  $\mathbf{p}$  accumulates the polarities of events that fall within its spatial region. The contrast of the IWE is measured by the gradient magnitude:

$$\|\nabla I_{iwe}\|^2 = \frac{1}{N_p} \sum_{i,j} (I_x^2(\mathbf{p}_{i,j}) + I_y^2(\mathbf{p}_{i,j})), \quad (4)$$

where  $\nabla I = (I_x, I_y)$  is the gradient of  $I_{iwe}$ ,  $N_p$  is the number of pixels,  $I_x \equiv \frac{\partial I}{\partial x}$ , and its magnitude is measured by the  $L2$  norm.

To stabilize convergence, we optimize zero-initialized residuals on top of the bilinearly up-sampled LR flows, constraining deviations from the initial estimates. The refined HR flow is obtained as

$$\mathbf{f}'_{hr} = \mathbf{U}(\mathbf{f}_{lr}) + \Delta\mathbf{f}_\theta, \quad (5)$$

where  $\Delta\mathbf{f}_\theta$  denotes the residual correction.

Combining equations 3, 4, 5, ICR executes per-scene optimization:

$$\arg \min_{\Delta\mathbf{f}_\theta} \|\nabla I_{iwe}(\mathbf{p}; \mathbf{f}'_{hr})\|^2. \quad (6)$$

Figure 3 demonstrates how CMax in HR space is able to recover fine-grained details. The increased spatial capacity from HR sensor simulation allows CMax to recover finer edge structures that are beyond the native sensor resolution.

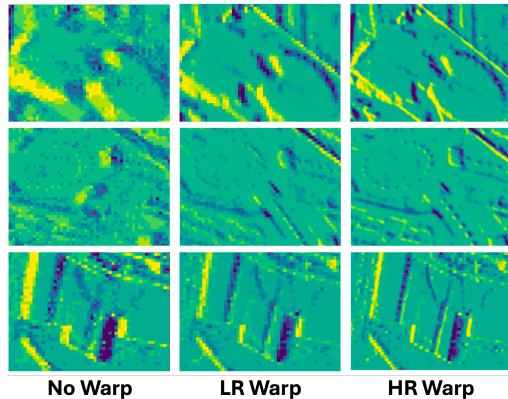


Figure 3: Comparison of IWEs under different warping regimes. No warp (left), LR CMax (middle), HR CMax via ICR (right).

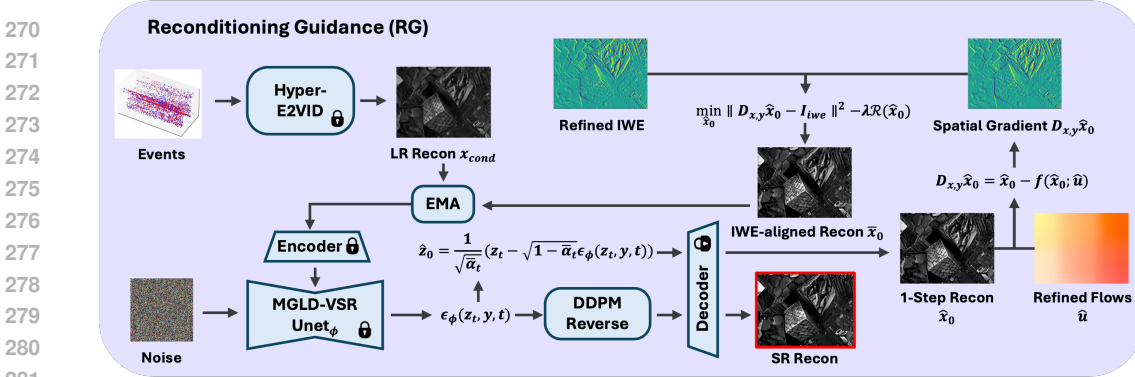


Figure 5: Our Reconditioning Guidance (RG) adds training-free control of the diffusion trajectory by carefully altering the conditioning image over diffusion steps. The 1-step reconstructions from intermediate latents are refined via the alignment of spatial gradients against the IWE, and the refined image is transferred to the conditioning image by exponential moving average.

However, the risk of event collapse in CMax increases substantially as the flow parameters under optimization acquire additional degrees of freedom. In the high-resolution domain, where the number of flow parameters grows by a factor of 16, collapse becomes inevitable. Fig. 4 (middle) illustrates this degeneracy: when CMax is applied directly in HR space, unconstrained flows concentrate events at regular spatial intervals, producing grid-like patterns in the IWE. Although these patterns yield high contrast scores, they show severely distorted scene structure and erase genuine details.

Our ICR is designed to mitigate event collapse in HR space via a local flow smoothing constraint that exploits the implicit smoothness of MLP to learn a continuous flow function:

$$u_{x,y}, v_{x,y} = G_{\theta}(\gamma_J(x), \gamma_J(y)), \quad (7)$$

where  $u, v$  are residual flow estimates at pixel  $(x, y)$ ,  $G_{\theta}$  is the MLP, and  $\gamma_J(\cdot)$  represents a positional encoding function with  $J$  frequency bases. We employ sinusoidal positional encodings Vaswani et al. (2023) on the query coordinates, retaining only low-frequency components to increase the similarity of positional encodings between local coordinates. We use  $J = 4$  in our ICR module and we also apply average pooling for final smoothing of flow residuals. Fig. 4 illustrates the benefits of ICR. ICR enables contrast maximization to be achieved in HR space by adopting a residual flow optimization process and providing implicit flow regularization to avoid event collapse while supporting fine-grained detail recovery from event data.

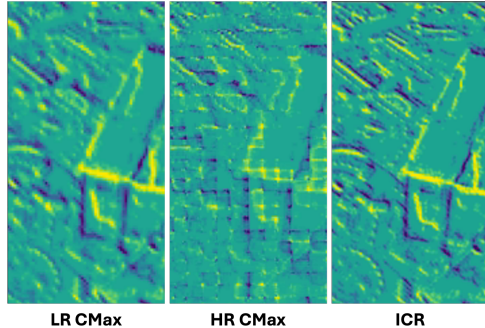


Figure 4: CMax in HR space degenerates with event collapse (middle). ICR uses frequency-constrained positional encodings in a lightweight MLP to regularize against degenerate HR flow solutions.

### 3.2 RECONDITIONING GUIDANCE

Reconditioning Guidance uses refined IWEs from ICR to guide the diffusion trajectory of a frozen MGLD-VSR Yang et al. (2024) model. Fig. 5 illustrates the RG process. At each diffusion step, noise predictions  $\epsilon_{\phi}(z_t, y, t)$  form one-step reconstructions  $\hat{z}_0$  which are decoded to image estimates  $\hat{x}_0$ . We then transfer fine-grained details from ICR by aligning  $\hat{x}_0$  to the refined IWE, yielding aligned frames  $\bar{x}_0$ .  $\bar{x}_0$  is then blended with conditioning image  $x_{cond}$  by exponential moving average (EMA) to update the conditioning  $y$  for the next diffusion step.

The alignment of  $\hat{x}_0$  to the refined IWE requires the matching of visual features across different modalities (intensity image vs IWE). To establish this link, we draw on Zhang et al. (2023), which formulates intensity image reconstruction as a linear inverse problem. Here, the IWE is approximated as the spatial derivative of the ground-truth intensity frame  $I_{gt}$ :

$$D_{x,y}I_{gt} \approx I_{IWE}(x, y), \quad (8)$$

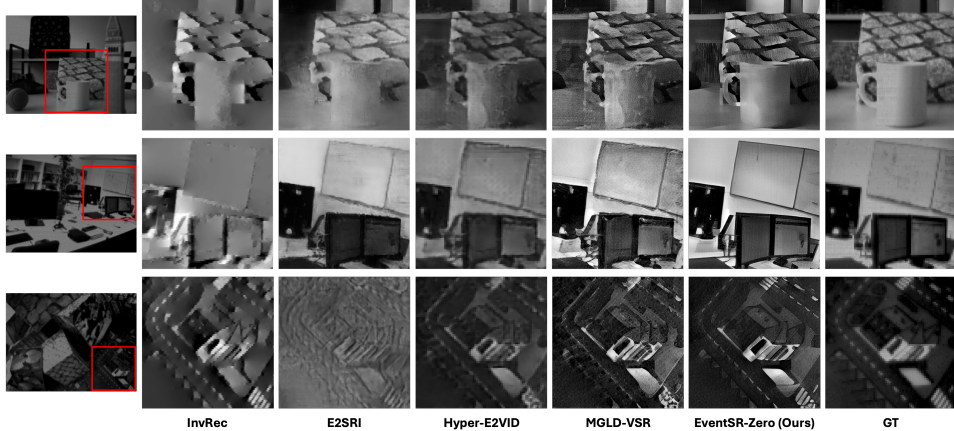


Figure 6: Qualitative Results on ECD Dataset.

where  $D_{x,y}$  denotes a finite difference operator. To align an image  $\ell$  with the IWE, we minimize:

$$\min_{\ell} \|D_{x,y}\ell - I_{iwe}\|^2 + \lambda\mathcal{R}(\ell), \quad (9)$$

where  $\ell$  is the image being optimized, and  $\lambda\mathcal{R}(\ell)$  is a regularization term. By enforcing flow-directed spatial gradient matching between  $\hat{x}_0$  and the refined IWE, fine-grained structures recovered by ICR are effectively transferred into the image estimate.

RG is applied at every diffusion step, where  $\hat{x}_0$  is refined for  $M = 100$  iterations using Eq. (9) with an L1 loss to limit deviations and a Total Variation loss Rudin et al. (1992) to suppress artifacts, producing detail-enhanced estimates  $\bar{x}_0$ .

The refined  $\bar{x}_0$  is then used as conditioning for the next diffusion step. While we experimented with self-guidance Epstein et al. (2023), its impact on SR outputs was negligible, likely because SR models rely heavily on the conditioning image, making them less responsive to intermediate latent adjustments. To address this limitation, our RG approach introduces a new control mechanism that directly adjusts the conditioning image during diffusion.

Our initial strategy was to fully replace the conditioning image with IWE-aligned reconstructions by setting  $\mathbf{x}_t^{\text{cond}} = \bar{x}_0$ . However, this approach caused errors from earlier diffusion steps to be reinjected into the conditioning, creating a feedback loop that amplified artifacts. To overcome this, we introduce an Exponential Moving Average (EMA) update that gradually incorporates IWE-aligned details while maintaining stability from past conditioning images. EMA effectively balances guidance and robustness by preventing abrupt changes that propagate errors across steps. Formally,

$$\mathbf{x}_{t-1}^{\text{cond}} = \eta\bar{x}_0 + (1 - \eta)\mathbf{x}_t^{\text{cond}}, \quad (10)$$

where  $\eta$  determines the trade-off between event guidance and stability. Large  $\eta$  values risk overemphasizing  $\bar{x}_0$  and amplifying artifacts, while small values underutilize IWE information.

Empirically, we found  $\eta = 0.05$  to be optimal for effective IWE guidance. The combination of spatial gradient alignment and EMA blending ensures that the diffusion process remains aligned with event-based observations over diffusion steps, guiding the VSR model to produce details that are structurally consistent with the scene.

## 4 EXPERIMENTS

In this section, we present both qualitative and quantitative results to highlight the effectiveness of the proposed EventSR-Zero model for  $4\times$  event video upscaling.

### 4.1 EXPERIMENT SETTINGS

We benchmark EventSR-Zero against three baseline models: 1) **InvRec**. The linear inverse reconstruction technique from Zhang et al. (2023). 2) **E2SRI**. The event-video super-resolution model

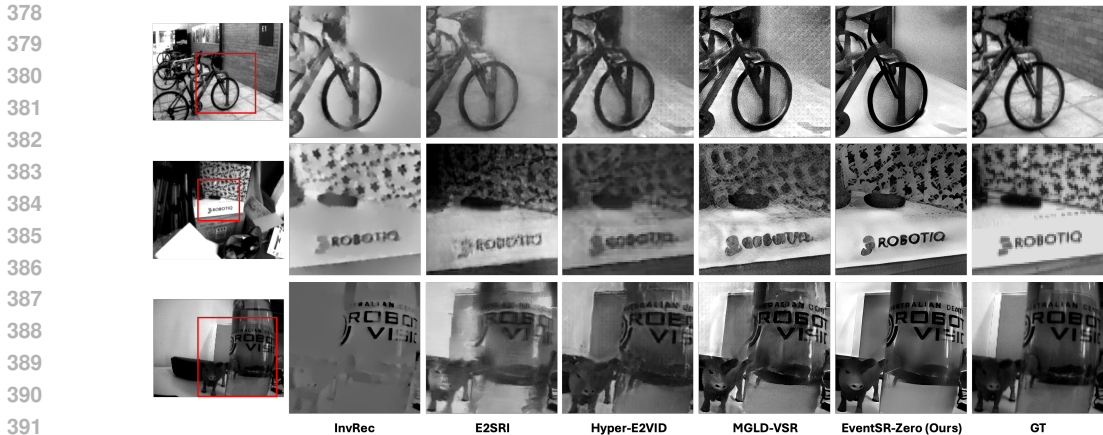


Figure 7: Qualitative Results on HQF Dataset.

introduced in Mostafavi et al. (2020). 3) **MGLD-VSR**. Our own baseline model with low-resolution frames generated by Hyper-E2VID Ercan et al. (2024) and upscaled using MGLD-VSR Yang et al. (2024).

#### 4.2 DATASET AND METRICS

For evaluation, we use sequences from two real-world datasets: the Event Camera Dataset (ECD) Mueggler et al. (2017) and the High-Quality Frames (HQF) Dataset Stoffregen et al. (2020b).

**Event Camera Dataset (ECD)**. ECD was captured using a DAVIS240C sensor to provide both events and frames at a resolution of  $240 \times 180$ , with ground truth frames recorded at 22 Hz. Following Rebecq et al. (2019), we select seven short sequences of static office scenes with 6-DOF camera motion.

**High-Quality Frames Dataset (HQF)**. HQF consists of indoor and outdoor scenes with diverse motion patterns, captured with a DAVIS240C sensor at  $240 \times 180$  resolution. Ground truth frames are recorded at 22.5 Hz and are carefully selected to ensure minimal motion blur.

**Metrics**. We evaluate our method using three full-reference metrics: 1) Mean Squared Error (MSE); 2) Structural Similarity (SSIM) Wang et al. (2004); 3) Learning Perceptual Image Patch Similarity (LPIPS) Zhang et al. (2018) to assess alignment with the true scenes structure. For the full-reference metrics, generated SR frames are downsampled to the original resolution for comparison with ground truth frames.

Since high-resolution ground truth frames are unavailable, we also apply two no-reference metrics: Naturalness Image Quality Evaluator (NIQE) Mittal et al. (2013) and BRISQUE Mittal et al. (2012) to assess the perceptual quality of SR images at high resolution. During evaluation, we perform histogram equalization on all reconstructions and ground truth images.

#### 4.3 QUALITATIVE RESULTS

The qualitative comparisons are presented in Fig. 6 and Fig. 7. GT refers to the LR ground truth frames. Our EventSR-Zero consistently reconstructs sharper edges and finer details compared to other methods. For example, our EventSR-Zero resolves the edges of the whiteboard more accurately in Fig. 6 row 2. This improvement is driven by our Reconditioning Guidance (RG), which effectively guides the diffusion process using high-resolution structures recovered by Implicit Contrast Refinement (ICR). This combination of fine-grained structure guidance and robust image priors of MGLD-VSR enables EventSR-Zero to produce clear and structured reconstructions of the whiteboard. The closeness of results from our EventSR-Zero to the structures in the ground truth (despite being low-resolution) highlights the effectiveness of ICR and RG in preserving alignment with scene structures derived from event data.

Table 1: Quantitative results of existing methods and our EventSR-Zero on sequences from ECD and HQF datasets. Best score is in bold and runner-up score is underlined.

	ECD					HQF				
	MSE ↓	SSIM ↑	LPIPS ↓	NIQE ↓	BRISQUE ↓	MSE ↓	SSIM ↑	LPIPS ↓	NIQE ↓	BRISQUE ↓
InvRec Zhang et al. (2023)	0.060	0.422	0.245	8.572	52.161	0.084	0.316	0.246	7.414	46.359
E2SRI Mostafavi et al. (2020)	0.066	0.419	0.226	6.236	28.131	0.073	0.313	0.231	6.379	33.004
MGLD-VSR Yang et al. (2024)	0.061	0.378	0.242	3.449	25.494	0.080	0.302	0.240	3.793	25.869
Ground Truth	-	-	-	9.312	68.742	-	-	-	9.075	65.852
EventSR-Zero (Ours)	<b>0.052</b>	<b>0.451</b>	<b>0.217</b>	<b>3.276</b>	<b>20.123</b>	<b>0.067</b>	<b>0.331</b>	<b>0.203</b>	<b>3.505</b>	<b>21.307</b>

The MGLD-VSR baseline in Fig. 6 column 4 further illustrates the impact of our ICR and RG modules. This baseline applies MGLD-VSR Yang et al. (2024) directly to Hyper-E2VID Ercan et al. (2024) outputs, and therefore is effectively our EventSR-Zero model without ICR and RG. We observe that on its own, MGLD-VSR tends to super-resolve existing artifacts and errors from Hyper-E2VID such as grid-like patterns as shown in Fig. 6 columns 3 and 4, without distinguishing these artifacts from true structures. In contrast, our EventSR-Zero with ICR and RG significantly sharpens edges and reduces artifacts to yield more refined SR outputs.

#### 4.4 QUANTITATIVE RESULTS

Quantitative results are shown in Tab. 1. Our EventSR-Zero outperforms competing baselines in MSE, SSIM, and LPIPS. Together, RG and ICR demonstrate complementary strengths in enhancing the alignment between super-resolved outputs and true scene structures. On no-reference metrics, our EventSR-Zero significantly outperforms other methods in NIQE Mittal et al. (2013) and BRISQUE Mittal et al. (2012) scores, indicating a substantial improvement in perceptual quality. The SR frames produced by EventSR-Zero exhibit enhanced sharpness and detail, contributing to a more realistic appearance.

**Discussion.** Since HR ground truth images are unavailable, full-reference evaluation necessitates downsampling the super-resolved (SR) outputs to match the lower-resolution ground truth for comparison. Although these low-resolution ground truth images retain the scene’s structural information, their poor visual quality makes them an imperfect reference for evaluating downsampled SR frames. This limitation is evident in their extremely high NIQE/BRISQUE scores (Table 1, Row 5).

#### 4.5 ABLATIONS

We ablate the performance contributions of RG, EMA and ICR in Tab. 2. The results show that both ICR and RG noticeably enhance reconstruction accuracy and visual quality across all metrics. ICR first recovers high-resolution scene structures

Table 2: Ablations of ICR and RG on HQF dataset.

RG	EMA	ICR	MSE ↓	SSIM ↑	LPIPS ↓	NIQE ↓	BRISQUE ↓
			0.080	0.302	0.240	3.793	25.869
✓	✓		0.072	0.324	0.215	3.652	23.070
✓		✓	0.093	0.266	0.291	3.591	21.847
✓	✓	✓	<b>0.067</b>	<b>0.331</b>	<b>0.203</b>	<b>3.505</b>	<b>21.307</b>

to offer more precise event guidance, followed by RG which leverages diffusion SR priors to produce high-quality event-guided details. Furthermore, EMA is essential for RG to produce accurate details with minimal hallucination. Without EMA (3rd row), the reference-based metrics (MSE/SSIM/LPIPS) worsen significantly due to the emergence of severe hallucinative artifacts.

## 5 CONCLUSION

In summary, EventSR-Zero introduces a training-free framework for event-to-video super-resolution that directly exploits the temporal richness of event data. The Implicit Contrast Refinement (ICR) module formulates a high-resolution contrast maximization space, regularized by a frequency-constrained implicit function, to recover high-frequency sub-pixel details while mitigating event collapse. Complementing this, the Reconditioning Guidance (RG) module steers the diffusion process by aligning the spatial gradients of intermediate latents with ICR outputs, ensuring that reconstructed frames preserve high-resolution scene structures. Together, ICR and RG enable structurally consistent, high-quality reconstructions from low-resolution events without requiring high-resolution training data, demonstrating the potential of training-free event-guided diffusion to bridge the gap between sparse event streams and high-resolution video in real-world vision applications.

## REFERENCES

- 486  
487  
488 A. Agrawal, R. Chellappa, and R. Raskar. An algebraic approach to surface reconstruction from  
489 gradient fields. In *Proceedings of the International Conference on Computer Vision (ICCV)*,  
490 volume 1, pp. 174–181, 2005.
- 491 M. Aharon, M. Elad, and A. Bruckstein. K-svd: An algorithm for designing overcomplete dictio-  
492 naries for sparse representation. *IEEE Transactions on Signal Processing*, 54(11):4311–4322,  
493 2006.
- 494 P. Bardow, A. J. Davison, and S. Leutenegger. Simultaneous optical flow and intensity estimation  
495 from an event camera. In *Proceedings of the IEEE Conference on Computer Vision and Pattern  
496 Recognition (CVPR)*, pp. 884–892, 2016.
- 497 S. Barua, Y. Miyatani, and A. Veeraraghavan. Direct face detection and video reconstruction from  
498 event cameras. In *Proceedings of the IEEE Winter Conference on Applications of Computer Vision  
499 (WACV)*, pp. 1–9, 2016.
- 500 Yuting Bi, Yan Zhang, Yebin Li, and Dong Chen. High dynamic range surveillance with an event  
501 camera. In *IEEE International Conference on Image Processing (ICIP)*, pp. 4350–4354, 2019.
- 502  
503 P. R. G. Cadena, Y. Qian, C. Wang, and M. Yang. Spade-e2vid: Spatially-adaptive denormalization  
504 for event-based video reconstruction. *IEEE Transactions on Image Processing*, 30:2488–2500,  
505 2021. doi: 10.1109/TIP.2021.3041899.
- 506 M. Cook, L. Gugelmann, F. Jug, C. Krautz, and A. Steger. Interacting maps for fast visual interpre-  
507 tation. In *Proceedings of the International Joint Conference on Neural Networks (IJCNN)*, pp.  
508 770–776, 2011.
- 509 Peiqi Duan, Zihao W. Wang, Xinyu Zhou, Yi Ma, and Boxin Shi. Eventzoom: Learning to de-  
510 noise and super resolve neuromorphic events. In *Proceedings of the IEEE/CVF Conference on  
511 Computer Vision and Pattern Recognition (CVPR)*, pp. 12824–12833, June 2021.
- 512  
513 Dave Epstein, Allan Jabri, Ben Poole, Alexei A. Efros, and Aleksander Holynski. Diffusion self-  
514 guidance for controllable image generation, 2023. URL [https://arxiv.org/abs/2306.](https://arxiv.org/abs/2306.00986)  
515 [00986](https://arxiv.org/abs/2306.00986).
- 516 Burak Ercan, Onur Eker, Canberk Saglam, Aykut Erdem, and Erkut Erdem. Hypere2vid: Improving  
517 event-based video reconstruction via hypernetworks. *IEEE Transactions on Image Processing*,  
518 33:1826–1837, 2024. ISSN 1941-0042. doi: 10.1109/tip.2024.3372460. URL [http://dx.](http://dx.doi.org/10.1109/TIP.2024.3372460)  
519 [doi.org/10.1109/TIP.2024.3372460](http://dx.doi.org/10.1109/TIP.2024.3372460).
- 520 Davide Falanga, Kevin Kleber, and Davide Scaramuzza. Dynamic obstacle avoidance for quadrotors  
521 with event cameras. *Science Robotics*, 5(40):eaaz9712, 2020.
- 522  
523 Benedek Forrai, Takahiro Miki, Daniel Gehrig, Marco Hutter, and Davide Scaramuzza. Event-  
524 based agile object catching with a quadrupedal robot, 2023. URL [https://arxiv.org/](https://arxiv.org/abs/2303.17479)  
525 [abs/2303.17479](https://arxiv.org/abs/2303.17479).
- 526 Andrew C. Freeman, Ketan Mayer-Patel, and Montek Singh. Accelerated event-based feature de-  
527 tection and compression for surveillance video systems. In *Proceedings of the 15th ACM Mul-  
528 timedia Systems Conference, MMSys ’24*, pp. 132–143, New York, NY, USA, 2024. Associa-  
529 tion for Computing Machinery. ISBN 9798400704123. doi: 10.1145/3625468.3647618. URL  
530 <https://doi.org/10.1145/3625468.3647618>.
- 531  
532 Ioannis Asmanis Kenneth Chaney Guillermo Gallego Kostas Daniilidis Friedhelm Hamann,  
533 Ziyun Wang. Motion-prior contrast maximization for dense continuous-time motion estimation.  
534 In *European Conference on Computer Vision (ECCV)*, 2024.
- 535 Guillermo Gallego, Henri Rebecq, and Davide Scaramuzza. A unifying contrast maximization  
536 framework for event cameras, with applications to motion, depth, and optical flow estimation.  
537 In *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 3867–3876.  
538 IEEE, June 2018. doi: 10.1109/cvpr.2018.00407. URL [http://dx.doi.org/10.1109/](http://dx.doi.org/10.1109/CVPR.2018.00407)  
539 [CVPR.2018.00407](http://dx.doi.org/10.1109/CVPR.2018.00407).

- 540 Yue Gao, Jiaxuan Lu, Siqi Li, Nan Ma, Shaoyi Du, Yipeng Li, and Qionghai Dai. Action recogni-  
541 tion and benchmark using event cameras. *IEEE Transactions on Pattern Analysis and Machine*  
542 *Intelligence*, 45(12):14081–14097, 2023. doi: 10.1109/TPAMI.2023.3300741.
- 543  
544 Yue Gao, Jiaxuan Lu, Siqi Li, Yipeng Li, and Shaoyi Du. Hypergraph-based multi-view action  
545 recognition using event cameras. *IEEE Transactions on Pattern Analysis and Machine Intelli-*  
546 *gence*, 46(10):6610–6622, 2024. doi: 10.1109/TPAMI.2024.3382117.
- 547 Mathias Gehrig and Davide Scaramuzza. Recurrent vision transformers for object detection with  
548 event cameras, 2023. URL <https://arxiv.org/abs/2212.05598>.
- 549  
550 Kunping Huang, Sen Zhang, Jing Zhang, and Dacheng Tao. Event-based simultaneous localiza-  
551 tion and mapping: A comprehensive survey, 2024. URL [https://arxiv.org/abs/2304.](https://arxiv.org/abs/2304.09793)  
552 09793.
- 553  
554 H. Kim, A. Handa, R. Benosman, S. Ieng, and A. Davison. Simultaneous mosaicing and tracking  
555 with an event camera. In *Proceedings of the British Machine Vision Conference (BMVC)*, 2014.
- 556  
557 Hanme Kim, Stefan Leutenegger, and Andrew J Davison. Real-time 3d reconstruction and 6-dof  
558 tracking with an event camera. In *Proceedings of the European Conference on Computer Vision*  
559 *(ECCV)*, pp. 349–364, 2016.
- 560  
561 Yijin Li, Han Zhou, Bangbang Yang, Ye Zhang, Zhaopeng Cui, Hujun Bao, and Guofeng Zhang.  
562 Graph-based asynchronous event processing for rapid object recognition. In *Proceedings of the*  
*IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 934–943, October 2021.
- 563  
564 Alexey Mitrokhin, Cornelia Fermüller, Chahatkumar Parameshwara, and Yiannis Aloimonos.  
565 Event-based moving object detection and tracking. In *IEEE/RSJ International Conference on*  
566 *Intelligent Robots and Systems (IROS)*, pp. 8173–8179, 2019.
- 567  
568 Anish Mittal, Anush Krishna Moorthy, and Alan Conrad Bovik. No-reference image quality assess-  
569 ment in the spatial domain. *IEEE Transactions on Image Processing*, 21(12):4695–4708, 2012.  
570 doi: 10.1109/TIP.2012.2214050.
- 571  
572 Anish Mittal, Rajiv Soundararajan, and Alan C. Bovik. Making a “completely blind” image qual-  
573 ity analyzer. *IEEE Signal Processing Letters*, 20(3):209–212, 2013. doi: 10.1109/LSP.2012.  
2227726.
- 574  
575 S. M. Mostafavi, J. Choi, and K.-J. Yoon. Learning to super resolve intensity images from events. In  
576 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*,  
pp. 2768–2776, 2020. doi: 10.1109/CVPR42600.2020.00283.
- 577  
578 Elias Mueggler, Henri Rebecq, Guillermo Gallego, Tobi Delbruck, and Davide Scaramuzza. The  
579 event-camera dataset and simulator: Event-based data for pose estimation, visual odometry, and  
580 slam. *The International Journal of Robotics Research*, 36(2):142–149, February 2017. ISSN  
581 1741-3176. doi: 10.1177/0278364917691115. URL [http://dx.doi.org/10.1177/](http://dx.doi.org/10.1177/0278364917691115)  
582 0278364917691115.
- 583  
584 G. Munda, C. Reinbacher, and T. Pock. Real-time intensity-image reconstruction for event cameras  
585 using manifold regularisation. *International Journal of Computer Vision*, 126(12):1381–1393,  
2018.
- 586  
587 F. Paredes-Valles and G. C. de Croon. Back to event basics: Self-supervised learning of image  
588 reconstruction for event cameras via photometric constancy. In *Proceedings of the IEEE/CVF*  
589 *Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 3446–3455, 2021. doi:  
590 10.1109/CVPR46437.2021.00342.
- 591  
592 Chiara Plizzari, Mirco Planamente, Gabriele Goletto, Marco Cannici, Emanuele Gusso, Matteo  
593 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)*, pp. 19935–19947, June 2022.

- 594 H. Rebecq, R. Ranftl, V. Koltun, and D. Scaramuzza. High speed and high dynamic range video  
595 with an event camera. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 41(6):  
596 1377–1391, 2019. doi: 10.1109/TPAMI.2018.2873577.
- 597  
598 Henri Rebecq, Daniel Gehrig, and Davide Scaramuzza. ESIM: an open event camera simulator.  
599 *Conf. on Robotics Learning (CoRL)*, October 2018.
- 600  
601 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-  
602 resolution image synthesis with latent diffusion models, 2022. URL [https://arxiv.org/  
603 abs/2112.10752](https://arxiv.org/abs/2112.10752).
- 604  
605 Leonid Rudin, Stanley Osher, and Emad Fatemi. Nonlinear total variation based noise removal al-  
606 gorithms. *Physica D: Nonlinear Phenomena*, 60:259–268, 11 1992. doi: 10.1016/0167-2789(92)  
607 90242-F.
- 608  
609 C. Scheerlinck, H. Rebecq, D. Gehrig, N. Barnes, R. Mahony, and D. Scaramuzza. Fast image  
610 reconstruction with an event camera. In *Proceedings of the Winter Conference on Applications of  
611 Computer Vision (WACV)*, pp. 156–163, 2020. doi: 10.1109/WACV45572.2020.9093119.
- 612  
613 Christoph Scheerlinck, Robert Mahony, and Timothy O’Shea. Continuous-time intensity estimation  
614 using event cameras. In *Asian Conference on Computer Vision (ACCV)*, pp. 308–324, 2018.
- 615  
616 Shintaro Shiba, Yoshimitsu Aoki, and Guillermo Gallego. Event collapse in contrast maximization  
617 frameworks. *Sensors*, 22(14):5190, July 2022a. ISSN 1424-8220. doi: 10.3390/s22145190. URL  
618 <http://dx.doi.org/10.3390/s22145190>.
- 619  
620 Shintaro Shiba, Yoshimitsu Aoki, and Guillermo Gallego. *Secrets of Event-Based Optical Flow*,  
621 pp. 628–645. Springer Nature Switzerland, 2022b. ISBN 9783031197970. doi: 10.1007/  
622 978-3-031-19797-0\_36. URL [http://dx.doi.org/10.1007/978-3-031-19797-0\\_  
623 36](http://dx.doi.org/10.1007/978-3-031-19797-0_36).
- 624  
625 T. Stoffregen, C. Scheerlinck, D. Scaramuzza, T. Drummond, N. Barnes, L. Kleeman, and R. Ma-  
626 hony. Reducing the sim-to-real gap for event cameras. In *Proceedings of the European Conference  
627 on Computer Vision (ECCV)*, 2020a.
- 628  
629 Timo Stoffregen, Cedric Scheerlinck, Davide Scaramuzza, Tom Drummond, Nick Barnes, Lindsay  
630 Kleeman, and Robert Mahony. Reducing the sim-to-real gap for event cameras. In *Computer  
631 Vision – ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Pro-  
632 ceedings, Part XXVII*, pp. 534–549, Berlin, Heidelberg, 2020b. Springer-Verlag. ISBN 978-3-  
633 030-58582-2. doi: 10.1007/978-3-030-58583-9\_32. URL [https://doi.org/10.1007/  
634 978-3-030-58583-9\\_32](https://doi.org/10.1007/978-3-030-58583-9_32).
- 635  
636 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez,  
637 Lukasz Kaiser, and Illia Polosukhin. Attention is all you need, 2023. URL [https://arxiv.  
638 org/abs/1706.03762](https://arxiv.org/abs/1706.03762).
- 639  
640 Aayush Atul Verma, Bharatesh Chakravarthi, Arpitsinh Vaghela, Hua Wei, and Yezhou Yang.  
641 etram: Event-based traffic monitoring dataset, 2024. URL [https://arxiv.org/abs/  
642 2403.19976](https://arxiv.org/abs/2403.19976).
- 643  
644 Jianyi Wang, Zongsheng Yue, Shangchen Zhou, Kelvin C.K. Chan, and Chen Change Loy. Exploit-  
645 ing diffusion prior for real-world image super-resolution. 2024.
- 646  
647 L. Wang, T.-K. Kim, and K.-J. Yoon. Eventsr: From asynchronous events to image reconstruc-  
648 tion, restoration, and super-resolution via end-to-end adversarial learning. In *Proceedings of the  
649 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 8315–8325,  
650 2020. doi: 10.1109/CVPR42600.2020.00835.
- 651  
652 Zhou Wang, A.C. Bovik, H.R. Sheikh, and E.P. Simoncelli. Image quality assessment: from error  
653 visibility to structural similarity. *IEEE Transactions on Image Processing*, 13(4):600–612, 2004.  
654 doi: 10.1109/TIP.2003.819861.

648 Wenming Weng, Yueyi Zhang, and Zhiwei Xiong. Event-based video reconstruction using trans-  
649 former. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*,  
650 2021.

651 Xi Yang, Chenhang He, Jianqi Ma, and Lei Zhang. Motion-guided latent diffusion for temporally  
652 consistent real-world video super-resolution. 2024.

653 Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable  
654 effectiveness of deep features as a perceptual metric. In *CVPR*, 2018.

655 Zelin Zhang, Anthony J. Yezzi, and Guillermo Gallego. Formulating event-based image reconstruc-  
656 tion as a linear inverse problem with deep regularization using optical flow. *IEEE Trans. Pattern*  
657 *Anal. Mach. Intell.*, 45(7):8372–8389, July 2023. ISSN 0162-8828. doi: 10.1109/TPAMI.2022.  
658 3230727. URL <https://doi.org/10.1109/TPAMI.2022.3230727>.

659 Shangchen Zhou, Peiqing Yang, Jianyi Wang, Yihang Luo, and Chen Change Loy. Upscale-A-  
660 Video: Temporal-consistent diffusion model for real-world video super-resolution. In *CVPR*,  
661 2024.

662 Nikola Zubić, Mathias Gehrig, and Davide Scaramuzza. State space models for event cameras,  
663 2024. URL <https://arxiv.org/abs/2402.15584>.

664  
665  
666  
667  
668  
669  
670  
671  
672  
673  
674  
675  
676  
677  
678  
679  
680  
681  
682  
683  
684  
685  
686  
687  
688  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
700  
701