

UNIINR: UNIFYING SPATIAL-TEMPORAL INR FOR RS VIDEO CORRECTION, DEBLUR, AND INTERPOLATION WITH AN EVENT CAMERA

Anonymous authors

Paper under double-blind review

ABSTRACT

1 Images captured by rolling shutter (RS) cameras under fast camera motion often
 2 contain obvious image distortions and blur, which can be modeled as a row-wise
 3 combination of a sequence of global shutter (GS) frames within the exposure time.
 4 Naturally, recovering high-frame-rate GS sharp frames from an RS blur image
 5 needs to simultaneously consider RS correction, deblur, and frame interpolation.
 6 Tackling this task is nontrivial, and to the best of our knowledge, no feasible solu-
 7 tions exist by far. A naive way is to decompose the whole process into separate tasks
 8 and simply cascade existing methods; however, this results in cumulative errors
 9 and noticeable artifacts. Event cameras enjoy many advantages, *e.g.*, high temporal
 10 resolution, making them potential for our problem. To this end, we propose the
 11 **first** and novel approach, named **UniINR**, to recover arbitrary frame-rate sharp
 12 GS frames from an RS blur image and paired event data. Our key idea is *unifying*
 13 *spatial-temporal implicit neural representation (INR) to directly map the position*
 14 *and time coordinates to RGB values to address the interlocking degradations in*
 15 *the image restoration process*. Specifically, we introduce spatial-temporal implicit
 16 encoding (STE) to convert an RS blur image and events into a spatial-temporal
 17 representation (STR). To query a specific sharp frame (GS or RS), we embed
 18 the exposure time into STR and decode the embedded features pixel-by-pixel to
 19 recover a sharp frame. Our method features a lightweight model with only 0.379M
 20 parameters, and it also enjoys high inference efficiency, achieving 2.83ms/frame
 21 in $31 \times$ frame interpolation of an RS blur frame. Extensive experiments show that
 22 our method significantly outperforms prior methods.

23 1 INTRODUCTION

24 Most consumer-level cameras based on CMOS sensors rely on a rolling shutter (RS) mechanism.
 25 These cameras dominate the market owing to their benefits, *e.g.*, low power consumption (Janesick
 26 et al., 2009). In contrast to the global shutter (GS) cameras, RS cameras capture pixels row by row;
 27 thus, the captured images often suffer from obvious spatial distortions (*e.g.*, stretch) and blur under
 28 fast camera/scene motion. It has been shown that naively neglecting the RS effect often hampers the
 29 performance in many real-world applications (Hedborg et al., 2012; Lao & Ait-Aider, 2020; Zhong
 30 et al., 2021; Zhou et al., 2022). In theory, an RS image can be formulated as a row-wise combination
 31 of sequential GS frames within the exposure time (Fan & Dai, 2021; Fan et al., 2023).

32 In this regard, it is meaningful to *recover high-frame-rate sharp GS frames from a single RS blur*
 33 *image* as the restored high-frame-rate sharp GS frames can directly facilitate many downstream tasks
 34 in practice. Intuitively, achieving this goal often requires simultaneously considering RS correction,
 35 deblurring, and frame interpolation. However, tackling this task is nontrivial because multiple
 36 degradations, such as RS distortion, motion blur, and temporal discontinuity (Meilland et al., 2013;
 37 Su & Heidrich, 2015), often co-exist for CMOS cameras (Zhong et al., 2021). The co-existence of
 38 various image degradations complicates the whole GS frame restoration process. ***To the best of our***
 39 ***knowledge, no practical solutions exist in the literature to date.*** A naive way is to decompose the
 40 whole process as separate tasks and simply cascading existing image enhancement networks can
 41 result in cumulative errors and noticeable artifacts. For example, a simple consideration of cascading
 42 a frame interpolation network (Bao et al., 2019) with an RS correction network produces degraded
 43 results, as previously verified in (Naor et al., 2022).

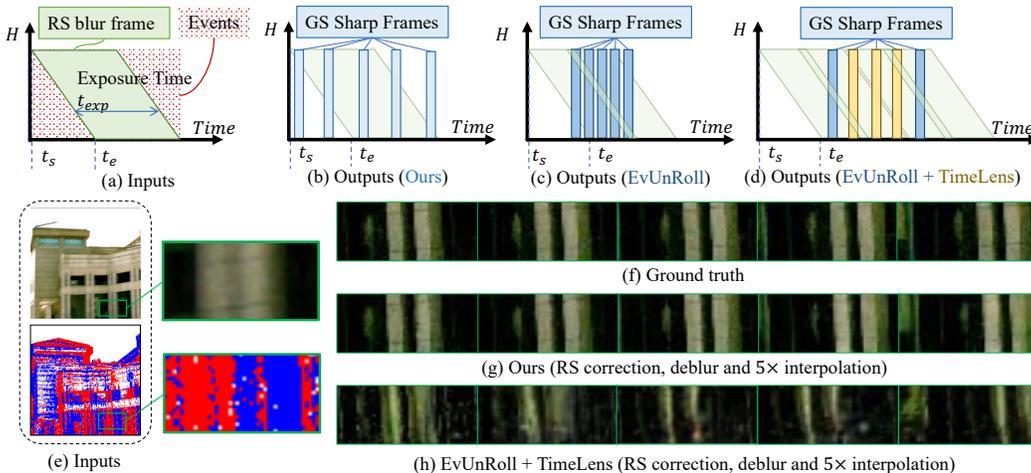


Figure 1: Inputs and the outputs of our method, EvUnRoll, and EvUnRoll+TimeLens. Inputs are shown in (a), which includes an RS blur image and events. t_s and t_e are the start and end timestamps of RS, and t_{exp} is the exposure time. Our outputs are shown in (b), which is a sequence of GS sharp frames during the whole exposure time of the RS blur image. (c) shows outputs of EvUnRoll, which can only recover the GS sharp frames in a limited time instead of the whole exposure time of the RS blur frame. (d) shows outputs of EvUnRoll+TimeLens. *More details are in Sec. C7 in Supp. Mat.*

44 Event cameras offer several advantages, such as high-temporal resolution, which make them suitable
 45 for various image restoration tasks (Wang et al., 2020; Zhou et al., 2022; Tulyakov et al., 2021; Song
 46 et al., 2023; 2022). eSL-Net (Wang et al., 2020) proposes an event-guided sparse learning framework
 47 to simultaneously achieve image super-resolution, denoising, and deblurring. TimeLens (Tulyakov
 48 et al., 2021) integrates a synthesis-based branch with a warp-based branch to boost the performance
 49 of the video frame interpolation. DeblurSR (Song et al., 2023) and E-CIR (Song et al., 2022) take ad-
 50 vantage of the high temporal resolution of events by converting a blurry frame into a time-to-intensity
 51 function, using spike representation and Lagrange polynomials, respectively. EvUnRoll (Zhou et al.,
 52 2022) leverages events as guidance to enhance RS correction by accounting for nonlinear motion
 53 during the desired timestamp. **However, these methods focus on either deburring or RS correction
 54 and can not recover arbitrary frame-rate sharp GS frames from a single RS blur image.** An
 55 example is depicted in Fig. 1 (h), showing that simply cascading event-guided RS correction model
 56 (e.g., EvUnroll (Zhou et al., 2022)) and interpolation model (e.g., TimeLens (Tulyakov et al., 2021))
 57 to recover high-frame-rate sharp GS frames results in obvious artifacts.

58 In this paper, we make the **first** attempt to propose a novel yet efficient learning framework, dubbed
 59 **UniINR**, that can **recover arbitrary frame-rate sharp GS frames from an RS blur image and events**.
 60 Our key idea is to learn a spatial-temporal *implicit neural representation (INR)* to *directly map the
 61 position and time coordinates to RGB values to address the co-existence of degradations in the image
 62 restoration process*. This makes it possible to exploit the spatial-temporal relationships from the
 63 inputs to achieve RS correction, deblur, and interpolation **simultaneously**. One distinct advantage
 64 of our method is that it is relatively lightweight with only **0.379M** parameters. We formulate
 65 the task —recovering high-frame-rate sharp GS frames from an RS blur image and paired event
 66 data —as a novel *estimation* problem, defined as a function, $F(x, t, \theta)$. Here, x denotes the pixel
 67 position (x, y) of an image, t denotes the timestamp during the exposure time, and θ denotes the
 68 function’s parameters. Our proposed framework consists of three parts: spatial-temporal implicit
 69 encoding (STE), exposure time embedding (ETE), and pixel-by-pixel decoding (PPD). Specifically,
 70 STE first utilizes sparse learning-based techniques (Wang et al., 2020) to extract a spatial-temporal
 71 representation (STR) θ from events and an RS blur image (Sec. 3.2.1). To query a specific sharp
 72 frame of RS or GS pattern, we then model the exposure information as a temporal tensor T in
 73 ETE (Sec. 3.2.2). Finally, PPD leverages an MLP to decode sharp frames from the STR and the
 74 temporal tensor T (Sec. 3.2.3), allowing for the generation of a sharp frame at any given exposure
 75 pattern (e.g., RS or GS). One notable advantage of our approach is its **high efficiency**, as it only
 76 requires using the STE once, regardless of the number of interpolation frames. In practice, as frame
 77 interpolation multiples rise from $1\times$ to $31\times$, the time taken increases from $31ms$ to $86ms$. Thus, at
 78 $31\times$ interpolation, each frame’s processing time is merely $2.8ms$, whereas the cascading approach
 79 (EvUnRoll + TimeLens) requires more than $177ms$ (Sec. 4.2).

80 We conduct a thorough evaluation of our proposed method, including both quantitative and qualitative
 81 analyses, using a higher resolution (256×256) dataset than that of the previous methods ($180 \times$
 82 240) (Song et al., 2023; 2022). Extensive experimental results demonstrate that our approach
 83 outperforms existing methods in RS correction, deblur, and interpolation (An example can be found
 84 in Fig. 1 (h)).

85 2 RELATED WORKS

86 2.1 EVENT-GUIDED IMAGE/VIDEO RESTORATION

87 **Event-guided Deblurring** Owing to the high temporal resolution afforded by events, prior stud-
 88 ies (Sun et al., 2022; Wang et al., 2020; Shang et al., 2021; Kim et al., 2022) have incorporated events
 89 into the task of deblurring. These works focus on the reconstruction of a single GS sharp frame
 90 from the GS blur frame, guided by event data. The work most analogous to ours is EvUnroll (Zhou
 91 et al., 2022), which first leverages event cameras for RS correction, leveraging their low latency
 92 benefits. Nonetheless, EvUnroll primarily focuses on RS correction, with its optional deblurring
 93 module equipped to handle minor motion blur and reconstruct a sharp frame at the midpoint of the
 94 exposure time.

95 **Event-guided Deblurring + Interpolation** These studies can be bifurcated based on the quantity
 96 of input GS blur frames: single GS frame (Xu et al., 2021; Song et al., 2022; 2023; Haoyu et al.,
 97 2020) or multiple GS frames (Pan et al., 2019; Zhang & Yu, 2022; Lin et al., 2020). The former,
 98 such as E-CIR (Song et al., 2022) and DeblurSR (Song et al., 2023), convert a GS blur frame into a
 99 time-to-intensity function while the latter, *e.g.*, EDI (Pan et al., 2019), LEDVDI (Lin et al., 2020),
 100 and EVDI (Zhang & Yu, 2022) are both built upon the event-based double integral model (Pan
 101 et al., 2019). However, these methods primarily target GS frames affected by motion blur, leading to
 102 performance degradation when dealing with spatially distorted and RS blur frames.

103 Recently, a contemporaneous study (Zhang et al., 2023) also focused on RS Correction, Deblur, and
 104 VFI. However, this research primarily concentrated on the individual performance of a single model
 105 across the three tasks, without extensive experimentation or investigation into handling all three tasks
 106 concurrently. This constitutes the most significant distinction from our method.

107 2.2 FRAME-BASED VIDEO RESTORATION FOR RS INPUTS

108 **RS Correction + Interpolation** RSSR (Fan & Dai, 2021; Fan et al., 2023) is the first work that
 109 generates multiple GS frames from two consecutive RS frames by introducing bi-directional undistor-
 110 tion flows. CVR (Fan et al., 2022) estimates two latent GS frames from two consecutive RS frames,
 111 followed by motion enhancement and contextual aggregation before generating final GS frames.

112 **RS Correction + Deblurring** JCD (Zhong et al., 2021) proposes the first pipeline that employs
 113 warping and deblurring branches to effectively address the RS distortion and motion blur. How-
 114 ever, JCD’s motion estimation module, built upon the assumption of linear motion derived from
 115 DeepUnrollNet (Liu et al., 2020), encounters a significant performance degradation in real-world
 116 scenarios involving non-linear motion (Zhou et al., 2022). To eliminate the dependence of motion
 117 estimation, (Wang et al., 2022b) proposes a method that turns the RS correction into a rectification
 118 problem, which allows all pixels to start exposure simultaneously and end exposure line by line.
 119 *Differently, our method can recover arbitrary GS sharp frames during the exposure time of RS blur*
 120 *frames without the assumption of linear motion.*

121 2.3 IMPLICIT NEURAL REPRESENTATION (INR)

122 INR (Wang et al., 2021; Sitzmann et al., 2020; Chen et al., 2021; 2022; Lu et al., 2023) is proposed
 123 for parameterized signals (images, video or audio) in the coordinate-based representation, inspiring
 124 some researchers to explore the potential of INR in low-level vision tasks. LIIF (Chen et al., 2021)
 125 represents images as high-dimensional tensors and allows for upsampling at any scale through
 126 interpolation and decoding, followed by VideoINR (Chen et al., 2022), which extends LIIF to videos,
 127 enabling temporal and spatial upsampling at any scale. EG-VSR (Lu et al., 2023) incorporates
 128 events into the learning of INR to achieve random-scale video super-resolution. *Differently, we*
 129 *propose STE to directly map the position and time coordinates to RGB values to address the co-*
 130 *existence of degradations in the image restoration process. Our STE makes it possible to exploit the*
 131 *spatial-temporal relationships from the inputs to achieve RS correction, deblur, and interpolation*
 132 *simultaneously.*

133 3 METHODOLOGY

134 3.1 PROBLEM DEFINITION AND ANALYSIS

135 We formulate the task —*recovering arbitrary frame-rate sharp GS frames from an RS blur image and*
 136 *paired event data*—as a novel estimation problem, defined as a function, $F(\mathbf{x}, t, \theta)$. Here, \mathbf{x} denotes

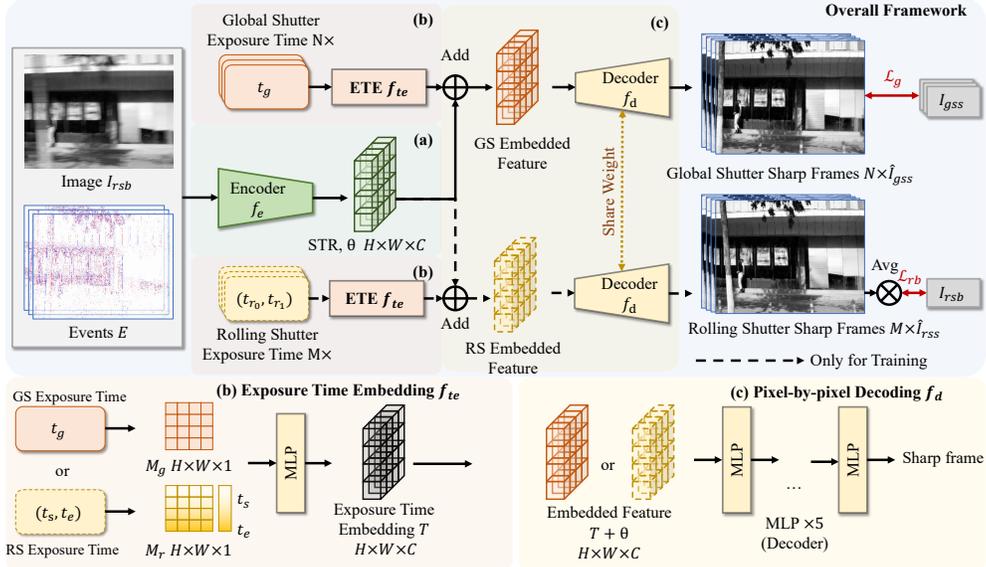


Figure 2: An overview of our framework, which consists of three parts, (a) the Spatial-Temporal Implicit Encoding (STE), (b) Exposure Time Embedding (ETE), and (c) Pixel-by-pixel decoding (PPD). Details of STE, ETE, and PPD are described in Sec. 3.2.1, Sec. 3.2.2, and Sec. 3.2.3. The inputs are an RS blur image $I_{r, sb}$ and events, and the outputs are a sequence of GS frames and RS frames. RS frames are predicted only in training.

137 the pixel position (x, y) of an image with a resolution of $H \times W$, t denotes the timestamp during the
 138 exposure time, and θ denotes the parameters. The intuition behind this formulation is that there exists
 139 a relationship between the RS blur/sharp frame and the GS blur/sharp frame. We now describe it.

140 By defining a function $F(\mathbf{x}, t, \theta)$ mapping the pixel position $\mathbf{x} = (x, y)$ and timestamp t to intensity
 141 or RGB value, we can obtain a GS sharp frame by inputting the desired timestamp \hat{t} during the
 142 exposure time to the function, which can be formulated as:

$$I_{g, \hat{t}} = F(\mathbf{x}, \hat{t}, \theta) \quad (1)$$

143 As an RS image can be formulated as a row-wise combination of sequential GS frames within the
 144 exposure time (Fan & Dai, 2021; Fan et al., 2023), we can assemble an RS sharp frame I_{r, t_s, t_e} from
 145 a sequence of GS sharp frames row by row given the RS start time t_s and the end time t_e . That is,
 146 the h -th row of an RS frame is the same as the h -th row of a GS frame at t_s^h , and the exposure start
 147 timestamp of the h -th row of an RS frame is $t_s^h = t_s + h \times (t_e - t_s)/H$. Therefore, we can formally
 148 describe an RS sharp frame as follows:

$$I_{r, t_s, t_e} = \left\{ F(\mathbf{x}, t_s^h, \theta) [h], h \in [0, H] \right\}. \quad (2)$$

149 In principle, a blur frame can be regarded as the temporal average of a sequence of sharp frames (Nah
 150 et al., 2017; Zhang et al., 2020). Thus, a GS blur frame $I_{g, t_g, t_{exp}}$, where t_g is the exposure start
 151 timestamp and t_{exp} is the exposure time, can be expressed as the average of a sequence of GS sharp
 152 frames during the exposure time t_{exp} , which can be formulated as:

$$I_{g, t_g, t_{exp}} = \frac{1}{t_{exp}} \int_{t_g}^{t_g + t_{exp}} F(\mathbf{x}, t, \theta) dt \approx \frac{1}{N} \sum_{i=0}^N I_{g, t_g + i \times t_{exp}/N}, \quad (3)$$

153 where N is the length of the GS frame sequence.

154 With the above formulation, an RS blur frame $I_{r, t_s \rightarrow t_e, t_{exp}}$ can thus be described based on the RS
 155 start time t_s , RS end time t_e , and exposure time of each scan line t_{exp} , as depicted in Fig. 1 (a).
 156 According to Eq. 2 and Eq. 3, the h -th row of an RS blur frame can be described as the temporal
 157 average of the h -th row in a sequence of GS sharp frames, which can be written as follows:

$$\begin{aligned} I_{r, t_s \rightarrow t_e, t_{exp}} &= \left\{ \frac{1}{t_{exp}} \int_{t_s^h}^{t_s^h + t_{exp}} F(\mathbf{x}, t, \theta) [h] dt, h \in [0, H] \right\} \\ &\approx \left\{ \frac{1}{N} \sum_{i=0}^N I_{g, t_s + i \times t_{exp}/N} [h], h \in [0, H] \right\}. \end{aligned} \quad (4)$$

158 An event stream E consists of a set of event $e = (x, y, t, p)$, where each event is triggered and
 159 recorded with the polarity p when the logarithmic brightness change at pixel (x, y) exceeds a certain

160 threshold C , which can be approximated as the differential of $F(\mathbf{x}, t, \theta)$ with respect to the time
 161 dimension. *For details about the principle of event cameras, refer to the Suppl. Mat.*

162 To use event data E as guidance, we need to address three challenges to estimate the mapping function
 163 $F(\mathbf{x}, t, \theta)$: **1**) how to find a function f_e to encode the input RS blur image and events to θ of the
 164 mapping function $F(\mathbf{x}, t, \theta)$; **2**) how to find a function f_{te} to represent the exposure information of
 165 desired RS or GS sharp frames as t of the mapping function $F(\mathbf{x}, t, \theta)$; **3**) how to find a function
 166 f_d to eliminate the need to input position information of desired RS or GS sharp frames as p of
 167 the mapping function $F(\mathbf{x}, t, \theta)$. Therefore, our goal is to estimate f_e , f_{te} , and f_d in order to get a
 168 mapped result, which can be formulated as:

$$I = F(\mathbf{x}, t, \theta) = F(\mathbf{x}, t, f_e(E, I_{rsb})) = F(\mathbf{x}, f_{te}(t), f_e(E, I_{rsb})) = f_d(f_{te}(t), f_e(E, I_{rsb})). \quad (5)$$

169 In the following section, we describe our framework based on Eq. 5 by substantiating f_e , f_{te} , and f_d .

170 3.2 PROPOSED FRAMEWORK

171 An overview of our UniINR framework is depicted in Fig. 2, which takes an RS blur image I_{rsb}
 172 and paired events E as inputs and outputs N sharp GS frames $\{I_{gss}\}_{i=0}^N$ with a high-frame-rate. To
 173 substantiate the defined functions f_e , f_{te} , and f_d , as mentioned in Sec. 3.1, our proposed framework
 174 consists of three components: **1**) Spatial-Temporal Implicit Encoding (STE), **2**) Exposure Time
 175 Embedding (ETE), and **3**) Pixel-by-pixel Decoding (PPD). Specifically, we first introduce an STE
 176 with deformable convolution (Wang et al., 2022a) to encode the RS blur frame and events into a
 177 spatial-temporal representation (STR) (Sec. 3.2.1). To provide exposure temporal information for
 178 STR, we embed the exposure start timestamp of each pixel from the GS or RS by ETE. (Sec. 3.2.2).
 179 Lastly, the PDD module adds ETE to STR to generate RS or GS sharp frames (Sec. 3.2.3). We now
 180 describe these components in detail.

181 3.2.1 SPATIAL-TEMPORAL IMPLICIT ENCODING (STE)

182 Based on the analysis in Sec. 3.1, we conclude that the RS blur frame I_{rsb} and events E collectively
 183 encompass the comprehensive spatial-temporal information during the exposure process. In this
 184 section, we aim to extract a spatial-temporal implicit representation θ that can effectively capture the
 185 spatial-temporal information from the RS blur frame I_{rsb} and events E .

186 To achieve this, we need to consider two key factors: (1) extracting features for the multi-task
 187 purpose and (2) estimating motion information. For the first factor, we draw inspiration from eSL-
 188 Net (Wang et al., 2020), which effectively utilizes events to simultaneously handle deblur, denoise,
 189 and super-resolution tasks. Accordingly, we design a sparse-learning-based backbone for the encoder.
 190 Regarding the second factor, previous works (Fan & Dai, 2021; Fan et al., 2022; 2023) commonly
 191 use optical flow for motion estimation in RS correction and interpolation tasks. However, optical
 192 flow estimation is computationally demanding (Gehrig et al., 2021; Zhu et al., 2019; Sun et al., 2018),
 193 making it challenging to incorporate it into the multiple task framework for RS cameras due to the
 194 complex degradation process. As an efficient alternative, we employ deformable convolution (Wang
 195 et al., 2022a) in our encoder to replace the optical flow estimation module. We adopt a 3D tensor with
 196 a shape of $H \times W \times C$ as the STR θ , which can effectively address the interlocking degradations
 197 encountered in the image restoration process with a sparse-learning-based backbone and deformable
 198 convolution, as formulated as $\theta = f_e(E, I_{rsb})$ in Eq. 5. *More details in the Suppl. Mat.*

199 3.2.2 EXPOSURE TIME EMBEDDING (ETE)

200 As depicted in Fig. 2 (b), the primary objective of the ETE module is to incorporate the exposure
 201 time of either a rolling shutter (RS) frame (t_s, t_e) or a global shutter (GS) frame (t_g) by employing
 202 an MLP layer, resulting in the generation of a temporal tensor T . To achieve this, we design an ETE
 203 module, denoted as f_{te} , which takes the GS exposure time t_g as input and produces the GS temporal
 204 tensor $T_g = f_{te}(t_g)$. Similarly, for RS frames, $T_r = f_{te}(t_{rs}, t_{re})$ represents the RS temporal tensor,
 205 which is only used in training. The process begins by converting the exposure process information
 206 into a timestamp map, with a shape of $H \times W \times 1$. Subsequently, the timestamp map is embedded
 207 by increasing its dimensionality to match the shape of the STR. This embedding procedure allows
 208 for the integration of the exposure time information into the STR representation. We now explain
 209 the construction of timestamp maps for both GS and RS frames and describe the embedding method
 210 employed in our approach.

211 **GS Timestamp Map:** In GS sharp frames, all pixels are exposed simultaneously, resulting in the
 212 same exposure timestamps for pixels in different positions. Given a GS exposure timestamp t_g , the
 213 GS timestamp map M_g can be represented as $M_g[h][w] = t_g$, where h and w denote the row and
 214 column indices, respectively.

215 **RS Timestamp Map:** According to the analysis in Sec. 3.1, pixels in RS frames are exposed line
 216 by line, and pixels in different rows have different exposure start timestamps. Given RS exposure
 217 information with start time t_s and RS end time t_e , the RS timestamp map can be represented as
 218 $M_r[h][w] = t_s + (t_e - t_s) \times h/H$, where h, w, H denote the row and column indices and height of
 219 the image, respectively.

220 **Time Embedding:** The timestamp maps, M_r and M_g , represent the timestamps of each pixel in
 221 a specific frame (RS or GS) with a shape of $H \times W \times 1$. However, the timestamp map is a high-
 222 frequency variable and can pose challenges for learning neural networks (Vaswani et al., 2017). Some
 223 approaches (Vaswani et al., 2017; Wang et al., 2021) propose a combination function of sine and
 224 cosine to encode the positional embedding. Nonetheless, calculating the derivative of the positional
 225 embedding is difficult, limiting its practical application to image enhancement tasks. In this paper,
 226 we utilize a one-layer MLP to increase the dimension for embedding. The whole embedding process
 227 is formulated as $T_g = f_{te}(t_g)$ for GS frames, and $T_r = f_{te}(t_{r_s}, t_{r_e})$ for RS frames, as depicted in
 228 Fig. 2(b). The MLP consists of a single layer that maps the timestamp map M_r or M_g to the same
 229 dimension $H \times W \times C$ as the spatial-temporal representation (STR) θ , as described in Sec. 3.2.1.

230 3.2.3 PIXEL-BY-PIXEL DECODING (PPD)

231 As shown in Fig. 2 (c), the goal of PPD is to efficiently query a sharp frame from STR θ by the
 232 temporal tensor T . It is important that the encoder is invoked only once for N times interpolation,
 233 while the decoder is called N times. Therefore, the efficiency of this query is crucial for the overall
 234 performance. The query’s inputs θ capture the global spatial-temporal information, and T captures
 235 the temporal information of the sharp frame (GS or RS). Inspired by previous works (Mildenhall
 236 et al., 2021; Chen et al., 2021), we directly incorporate the temporal tensor T into the STR θ to obtain
 237 an embedded feature with a shape of $H \times W \times C$ for each query. This additional embedded feature
 238 combines the global spatial-temporal information with the local exposure information, enabling
 239 straightforward decoding to obtain a sharp frame. To avoid the need for explicit positional queries,
 240 we employ a pixel-by-pixel decoder. The decoder, denoted as f_d in Eq. 5, employs a simple 5-layer
 241 MLP $f_{mlp}^{\odot 5}$ architecture. The reconstructed output I after decoding can be described in Eq. 6, where
 242 \oplus means element-wise addition.

$$I = f_d(f_{te}(t), f_e(E, I_{rsb})) = f_d(T, \theta) = f_{mlp}^{\odot 5}(T \oplus \theta). \quad (6)$$

243 3.2.4 LOSS FUNCTION

244 **RS Blur Image-guided Integral Loss:** Inspired by EVDI (Zhang & Yu, 2022), we formulate the
 245 relationship between RS blur frames and RS sharp frames. Given a sequence of RS sharp frames
 246 generated from the decoder, the input RS blur frame $I_{rsb} = \frac{1}{M} \sum_{i=1}^M (\hat{I}_{r_{ss}}^i)$, where M represents the
 247 length of the RS image sequence. In this way, we can formulate the blur frame guidance integral
 248 loss between the reconstructed RS blur frame and the original RS blur frame as $\mathcal{L}_b = \mathcal{L}_c(\hat{I}_{rsb}, I_{rsb})$,
 249 where \mathcal{L}_c indicates *Charbonnier loss* (Lai et al., 2018).

250 **Total Loss:** Apart from RS blur image-guided integral loss \mathcal{L}_b , we incorporate a reconstruction loss
 251 \mathcal{L}_{re} to supervise the reconstructed GS sharp frames. Our method consists of two losses: RS blur
 252 image-guided integral loss and the reconstruction loss, where λ_b, λ_{re} denote the weights of each loss:

$$\mathcal{L} = \lambda_b \mathcal{L}_b + \lambda_{re} \mathcal{L}_{re} = \lambda_b \mathcal{L}_c(\hat{I}_{rsb}, I_{rsb}) + \lambda_{re} \frac{1}{N} \sum_{k=1}^N \mathcal{L}_c(\hat{I}_{g_{ss}}^k, I_{g_{ss}}^k). \quad (7)$$

253 4 EXPERIMENTS

254 **Implementation Details:** We utilize the Adam optimizer (Kingma & Ba, 2014) for all experiments,
 255 with learning rates of $1e - 4$ for both Gev-RS (Zhou et al., 2022) and Fastec-RS (Liu et al., 2020)
 256 datasets. Using two NVIDIA RTX A5000 GPU cards, we train our framework across 400 epochs
 257 with a batch size of two. In addition, we use the mixed precision (Micikevicius et al., 2017) training
 258 tool provided by PyTorch (Paszke et al., 2017), which can speed up our training and reduce memory
 259 usage. PSNR and SSIM (Wang et al., 2004) are used to evaluate the reconstructed results.

260 **Datasets: 1) Gev-RS dataset (Zhou et al., 2022)** features GS videos at 1280×720 , 5700 fps. We
 261 reconstruct such frames and events from the original videos, downsampling to 260×346 (Scheerlinck
 262 et al., 2019). Events and RS blur frames are synthesized using vid2e (Gehrig et al., 2020). We adopt
 263 EvUnroll’s (Zhou et al., 2022) 20/9 train/test split. **2) Fastec-RS dataset (Liu et al., 2020)** offers GS
 264 videos at 640×480 , 2400 fps. We apply identical settings for resizing, event creation, and RS blur.
 265 The dataset is split into 56 training and 20 testing sequences. **3) Real-world dataset (Zhou et al.,
 266 2022)** is the only available real dataset, containing four videos with paired RS frames and events.
 267 Due to the lack of ground truth, it offers only quantitative visualizations. *More details in Suppl. Mat..*

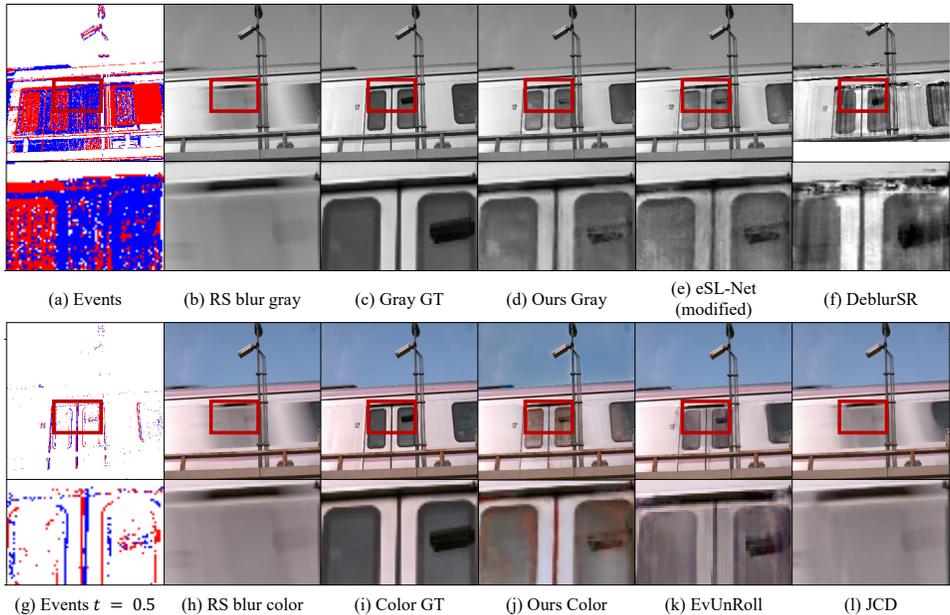


Figure 3: Visual Comparisons on RS correction and deblurring on Gev-RS (Zhou et al., 2022) dataset. The image resolution of DeblurSR (Song et al., 2023) is 180×240 .

268 4.1 COMPARISON WITH SOTA METHODS

269 Our experiments are conducted on both simulated and real datasets. While the simulated dataset
 270 enables us to obtain accurate quantitative results, evaluating on the real dataset offers insights into the
 271 generation ability of our method.

272 We compare our methods with recent methods with two different settings in these two datasets: **(I)**
 273 the experiment with a single GS sharp frame result, including JCD (Zhong et al., 2021) (frame-
 274 based RS correction and deblurring), EvUnroll (Zhou et al., 2022) (event-guided RS correction)
 275 and eSL-Net (Wang et al., 2020) (event-guided deblurring). **(II)** the experiment with a sequence of
 276 GS sharp frames result, which includes DeblurSR (Song et al., 2023) (event-guided deblurring and
 277 interpolation), and the combination of EvUnroll (Zhou et al., 2022) and TimeLens (Tulyakov et al.,
 278 2021) (event-guided video frame interpolation). In addition, we test our model’s generation ability by
 279 comparing it with EvUnRoll (Zhou et al., 2022) using real data. While this real data is solely reserved
 280 for testing, both our model and EvUnRoll are trained on the simulation dataset. *More explanations of*
 281 *setting (II) are in Supp. Mat..*

282 We evaluate JCD, EvUnroll, TimeLens, and DeblurSR with the released code. We modified eSL-
 283 Net by adjusting its parameterization initialization method and removing the up-sampling module,
 284 allowing it to be well trained on our datasets. The outputs of eSL-Net and DeblurSR are grayscale
 285 frames, and the outputs of JCD, EvUnroll, and the combination of EvUnroll and TimeLens are RGB
 286 frames. For fairness, our network is trained with the input of grayscale and RGB images, respectively.

287 The quantitative results for experiments generating a single GS sharp frame ($1 \times$) and those producing
 288 a sequence of GS sharp frames ($3 \times$, $5 \times$, $9 \times$) are presented in Tab. 1. In comparison to methods
 289 that yield a single GS sharp frame, our approach exhibits remarkable performance in both gray and
 290 RGB frames, surpassing the best-performing methods (eSL-Net (Wang et al., 2020) in gray and
 291 EvUnroll (Zhou et al., 2022) in RGB) by **1.48dB** and **4.17dB** on the Gev-RS (Zhou et al., 2022)
 292 dataset, respectively. In scenarios where a sequence of GS sharp frames is produced, our method
 293 attains optimal performance for both gray and RGB frames, achieving an increase of up to **13.47dB**
 294 and **8.49dB** compared to DeblurSR (Song et al., 2023) and EvUnroll (Zhou et al., 2022)+Time-
 295 Lens (Tulyakov et al., 2021) on the Gev-RS (Zhou et al., 2022) dataset, respectively. The substantial
 296 performance decline of DeblurSR (Song et al., 2023) can be ascribed to the interdependence between
 297 RS correction and deblur. The performance reduction of EvUnroll+TimeLens can be accounted for
 298 by the accumulation of errors arising from this cascading network, as shown in Fig. 1(h).

299 The qualitative results, as depicted in Fig. 11, showcase the effectiveness of our proposed method
 300 on both grayscale and RGB inputs. These results demonstrate our approach’s ability to generate

Table 1: Quantitative results for RS correction, deblurring, and frame interpolation. TL refers to TimeLens Tulyakov et al. (2021). EU refers to EvUnroll Zhou et al. (2022). eSL-Net* represents a modified model based on eSL-Net Wang et al. (2020).

	Methods	Inputs		Params(M) ↓	Gev-RS		Fastec-RS	
		Frame	Event		PSNR ↑	SSIM ↑	PSNR ↑	SSIM ↑
1 ×	eSL-Net*	1 gray	✓	0.1360	31.64	0.9614	32.45	0.9186
	UniINR (Ours)	1 gray	✓	0.3790	33.12	0.9881	34.62	0.9390
	JCD	3 color	✗	7.1659	18.59	0.5781	21.31	0.6150
	EU	1 color	✓	20.83	26.18	0.8606	29.76	0.8693
	UniINR (Ours)	1 color	✓	0.3792	30.35	0.9714	33.64	0.9299
3 ×	DeblurSR	1 gray	✓	21.2954	17.64	0.554	21.17	0.5816
	UniINR (Ours)	1 gray	✓	0.3790	31.11	0.9738	33.23	0.9210
	EU + TL	2 color	✓	93.03	21.86	0.7057	24.81	0.7179
	UniINR (Ours)	1 color	✓	0.3792	28.36	0.9348	32.72	0.9147
5 ×	DeblurSR	1 gray	✓	21.2954	18.35	0.6107	22.86	0.6562
	UniINR (Ours)	1 gray	✓	0.3790	30.84	0.9673	32.82	0.9147
	EU + TL	2 color	✓	93.03	21.59	0.6964	24.46	0.7140
	UniINR (Ours)	1 color	✓	0.3792	28.41	0.9062	32.13	0.9053
9 ×	DeblurSR	1 gray	✓	21.2954	18.86	0.6502	23.96	0.7049
	UniINR (Ours)	1 gray	✓	0.3790	30.54	0.9579	32.21	0.9051
	EU + TL	2 color	✓	93.03	21.24	0.6869	23.99	0.7029
	UniINR (Ours)	1 color	✓	0.3792	27.21	0.8869	29.31	0.8590

Table 2: Quantitative comparison in PSNR, SSIM, and LPIPS on EvUnRoll simulation dataset (Zhou et al., 2022). The numerical results of DSUN, JCD, and EvUnRoll are provided by (Zhou et al., 2022).

Method	Frames	Event	Params(M) ↓	PSNR ↑	SSIM ↑	LPIPS ↓
DSUN (Liu et al., 2020)	2	✗	3.91	23.10	0.70	0.166
JCD (Zhong et al., 2021)	3	✗	7.16	24.90	0.82	0.105
EvUnRoll (Zhou et al., 2022)	1	✓	20.83	30.14	0.91	0.061
UniINR(Ours)	1	✓	0.38	30.61	0.9285	0.048

301 sharp frames devoid of RS distortion, yielding the most visually pleasing outcomes in challenging
 302 scenarios involving a fast-moving train with motion blur and RS distortion. Comparatively, the results
 303 of eSL-Net and EvUnroll exhibit discernible noise, particularly evident around the train door within
 304 the red region of Fig. 11. Another approach, JCD, falls short in recovering sharp frames within
 305 such complex scenes. This failure can be attributed to the insufficient availability of frame-based
 306 methods which rely on the assumption of linear motion. Furthermore, the results obtained using
 307 DebluSR (Song et al., 2023) display noticeable artifacts, particularly in the context of the moving
 308 train. These artifacts hinder satisfactory frame reconstruction in such dynamic environments.

309 *Bad case analysis:* The color distortion in Fig. 11 (j) can be attributed to the insufficient color
 310 information in the challenging scene of a fast-moving train. From the input (Fig. 11 (h)), it can be
 311 noticed that the degree of motion blur is extremely severe and the blurry frame cannot provide valid
 312 color information. Furthermore, according to the principle of the generation of the event, the event is
 313 triggered by intensity change and it cannot provide color information.

314 **EvUnRoll Simulation Dataset:** To achieve a more equitable comparison with EvUnRoll, we evaluate
 315 our method on the simulated dataset employed by EvUnRoll, shown in Tab. 2. It’s important to
 316 emphasize that the dataset includes paired data consisting of RS blur, RS sharp, and GS sharp. For our
 317 model’s training, we specifically utilize the paired images of RS blur and GS sharp. As a one-stage
 318 approach, our method directly transforms an RS-blurred image into a GS-sharp image avoiding
 319 accumulated error, and thus has better performance.

320 **Real-world Dataset:** Fig. 4 shows real-world results. The input frame exhibits rolling shutter
 321 distortions, such as curved palette edges. In contrast, events show global shutter traits. Both
 322 our method and EvUnRoll correct these distortions effectively. Due to the lack of ground truth,
 323 quantitative analysis is not possible. Notably, our method avoids artifacts and errors, outperforming
 324 EvUnRoll in palette scenarios. *For further discussion please refer to the Supp. Mat..*

325 4.2 ABLATION AND ANALYTICAL STUDIES

326 **Importance of Exposure Time Embedding:** We conduct the experiments to evaluate the impact
 327 of learning-based position embedding, with a comparative analysis to sinusoid position embed-

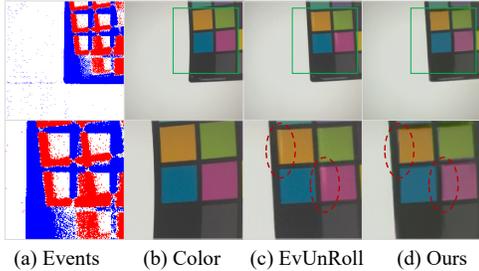


Figure 4: Visualization results in a real-world dataset (Zhou et al., 2022). (a) is the events visualization results. (b) are the input RGB images that have clear rolling shutter distortions. (c) is the output of EvUnRoll. (d) are the outputs of our method. The red circle in (c) has color distortion.

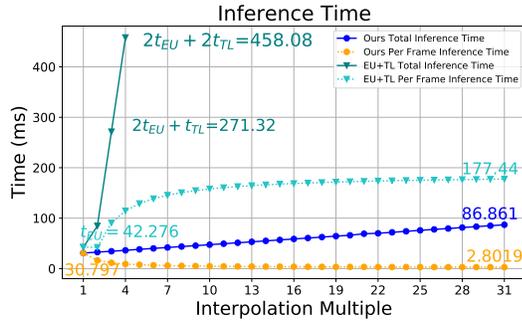


Figure 5: Comparison of inference time of our method with EvUnroll + TimeLens. t_{EU} and t_{TL} represent the respective inference times of EvUnRoll and TimeLens. The axes represent frame interpolation multiples ($1\times$ to $31\times$) and time. $2t_{EU}$ and $2t_{TL}$ means calling EvUnRoll twice and TimeLens twice.

328 ding (Vaswani et al., 2017). As indicated in Tab. 3, learning-based position embedding outperforms
 329 sinusoid position embedding, with advancements of up to **1.11dB** on average. This superior efficacy
 330 is attributable to the intrinsic adaptability of the learning-based position embedding.

331 **Importance of RS Blur Image-guided Integral Loss:** The effectiveness of the RS blur image-guided
 332 integral Loss across diverse interpolation settings is depicted in Tab. 4. The findings point towards
 333 the enhancement in PSNR for high interpolation configurations (e.g., $9\times$) upon employing this loss.

334 **Inference Speed:** Fig. 5 shows our method’s inference time across $1\times$ to $31\times$ interpolation. The
 335 total time rises modestly, e.g., from 30.8 ms at $1\times$ to 86.9 ms at $31\times$, a 2.8-fold increase for a
 336 31-fold interpolation. The average frame time even decreases at higher multiples, reaching 2.8 ms at
 337 $31\times$. Compared to EvUnRoll (Zhou et al., 2022) and TimeLens (Tulyakov et al., 2021), our method
 338 is more computationally efficient, requiring only 72% of EvUnRoll’s 42.3 ms for RS correction and
 339 deblurring. For N -fold frame insertion using EvUnRoll + TimeLens, EvUnRoll is counted twice,
 340 and TimeLens $N - 2$ times. This advantage is amplified in high-magnification scenarios, where
 341 TimeLens costs 186.76 ms per call. Our calculations focus on GPU time, excluding data I/O, which
 342 further increases EvUnRoll and TimeLens’ time consumption. *More discussions are in Supp. Mat..*

343 Table 3: Ablation for learning-based position embedding.

	Position Embedding	PSNR	SSIM
1x	Sinusoid	32.46	0.9851
	Learning	33.12	0.9881
3x	Sinusoid	30.83	0.9723
	Learning	31.11	0.9738
5x	Sinusoid	30.70	0.9678
	Learning	30.84	0.9673
9x	Sinusoid	30.51	0.9560
	Learning	30.54	0.9579
		+1.11	+0.0059

Table 4: Ablation for the loss function.

	\mathcal{L}_b	PSNR	SSIM
1x	✗	33.12	0.9881
	✓	33.14	0.9844
3x	✗	31.11	0.9738
	✓	31.09	0.9768
5x	✗	30.84	0.9673
	✓	30.83	0.9784
9x	✗	30.54	0.9579
	✓	30.61	0.9538
		+0.060	+0.0063

344 **5 CONCLUSION**

345 This paper presented a novel approach that simultaneously uses events to guide rolling shutter frame
 346 correction, deblur, and interpolation. Unlike previous network structures that can only address one
 347 or two image enhancement tasks, our method incorporated all three tasks concurrently, providing
 348 potential for future expansion into areas such as image and video super-resolution and denoising.
 349 Furthermore, our approach demonstrated high efficiency in computational complexity and model size.
 350 Regardless of the number of frames involved in interpolation, our method only requires a single call
 351 to the encoder, and the model size is a mere 0.379M.

352 **Limitations** Our analysis utilizes simulated data and real-world datasets, the latter of which lacks
 353 ground truth. Acquiring real data with ground truth is challenging. In future work, we aim to address
 354 this limitation by employing optical instruments, such as spectroscopes, to obtain real-world data
 355 with ground truth for quantitative evaluation.

356 REFERENCES

- 357 Wenbo Bao, Wei-Sheng Lai, Chao Ma, Xiaoyun Zhang, Zhiyong Gao, and Ming-Hsuan Yang. Depth-
358 aware video frame interpolation. In *Proceedings of the IEEE/CVF Conference on Computer Vision*
359 *and Pattern Recognition*, pp. 3703–3712, 2019. 1
- 360 Yinbo Chen, Sifei Liu, and Xiaolong Wang. Learning continuous image representation with local
361 implicit image function. In *Proceedings of the IEEE/CVF conference on computer vision and*
362 *pattern recognition*, pp. 8628–8638, 2021. 3, 6
- 363 Zeyuan Chen, Yinbo Chen, Jingwen Liu, Xingqian Xu, Vidit Goel, Zhangyang Wang, Humphrey
364 Shi, and Xiaolong Wang. Videoinr: Learning video implicit neural representation for continuous
365 space-time super-resolution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and*
366 *Pattern Recognition*, pp. 2047–2057, 2022. 3, 18
- 367 Bin Fan and Yuchao Dai. Inverting a rolling shutter camera: bring rolling shutter images to high
368 framerate global shutter video. In *Proceedings of the IEEE/CVF International Conference on*
369 *Computer Vision*, pp. 4228–4237, 2021. 1, 3, 4, 5
- 370 Bin Fan, Yuchao Dai, Zhiyuan Zhang, Qi Liu, and Mingyi He. Context-aware video reconstruction
371 for rolling shutter cameras. In *Proceedings of the IEEE/CVF Conference on Computer Vision and*
372 *Pattern Recognition*, pp. 17572–17582, 2022. 3, 5
- 373 Bin Fan, Yuchao Dai, and Hongdong Li. Rolling shutter inversion: Bring rolling shutter images to
374 high framerate global shutter video. *IEEE Transactions on Pattern Analysis & Machine Intelligence*,
375 45(05):6214–6230, 2023. 1, 3, 4, 5
- 376 Daniel Gehrig, Mathias Gehrig, Javier Hidalgo-Carrió, and Davide Scaramuzza. Video to events:
377 Recycling video datasets for event cameras. In *Proceedings of the IEEE/CVF Conference on*
378 *Computer Vision and Pattern Recognition*, pp. 3586–3595, 2020. 6, 14
- 379 Mathias Gehrig, Mario Millhäusler, Daniel Gehrig, and Davide Scaramuzza. E-raft: Dense optical
380 flow from event cameras. In *2021 International Conference on 3D Vision (3DV)*, pp. 197–206.
381 IEEE, 2021. 5
- 382 Chen Haoyu, Teng Minggui, Shi Boxin, Wang Yizhou, and Huang Tiejun. Learning to deblur and
383 generate high frame rate video with an event camera. *arXiv preprint arXiv:2003.00847*, 2020. 3
- 384 Johan Hedborg, Per-Erik Forssén, Michael Felsberg, and Erik Ringaby. Rolling shutter bundle
385 adjustment. In *2012 IEEE Conference on Computer Vision and Pattern Recognition*, pp. 1434–
386 1441. IEEE, 2012. 1
- 387 James Janesick, Jeff Pinter, Robert Potter, Tom Elliott, James Andrews, John Tower, John Cheng, and
388 Jeanne Bishop. Fundamental performance differences between cmos and ccd imagers: part iii. In
389 *Astronomical and Space Optical Systems*, volume 7439, pp. 47–72. SPIE, 2009. 1
- 390 Taewoo Kim, Jeongmin Lee, Lin Wang, and Kuk-Jin Yoon. Event-guided deblurring of unknown
391 exposure time videos. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv,*
392 *Israel, October 23–27, 2022, Proceedings, Part XVIII*, pp. 519–538. Springer, 2022. 3
- 393 Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint*
394 *arXiv:1412.6980*, 2014. 6
- 395 Wei-Sheng Lai, Jia-Bin Huang, Narendra Ahuja, and Ming-Hsuan Yang. Fast and accurate image
396 super-resolution with deep laplacian pyramid networks. *IEEE transactions on pattern analysis and*
397 *machine intelligence*, 41(11):2599–2613, 2018. 6
- 398 Yizhen Lao and Omar Ait-Aider. Rolling shutter homography and its applications. *IEEE transactions*
399 *on pattern analysis and machine intelligence*, 43(8):2780–2793, 2020. 1
- 400 Songnan Lin, Jiawei Zhang, Jinshan Pan, Zhe Jiang, Dongqing Zou, Yongtian Wang, Jing Chen, and
401 Jimmy Ren. Learning event-driven video deblurring and interpolation. In *Computer Vision–ECCV*
402 *2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part VIII 16*,
403 pp. 695–710. Springer, 2020. 3

- 404 Peidong Liu, Zhaopeng Cui, Viktor Larsson, and Marc Pollefeys. Deep shutter unrolling network.
405 In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp.
406 5941–5949, 2020. 3, 6, 8, 14
- 407 Yunfan Lu, Zipeng Wang, Minjie Liu, Hongjian Wang, and Lin Wang. Learning spatial-
408 temporal implicit neural representations for event-guided video super-resolution. *arXiv preprint*
409 *arXiv:2303.13767*, 2023. 3
- 410 Maxime Meilland, Tom Drummond, and Andrew I Comport. A unified rolling shutter and motion
411 blur model for 3d visual registration. In *Proceedings of the IEEE International Conference on*
412 *Computer Vision*, pp. 2016–2023, 2013. 1
- 413 Paulius Micikevicius, Sharan Narang, Jonah Alben, Gregory Diamos, Erich Elsen, David Garcia,
414 Boris Ginsburg, Michael Houston, Oleksii Kuchaiev, Ganesh Venkatesh, et al. Mixed precision
415 training. *arXiv preprint arXiv:1710.03740*, 2017. 6
- 416 Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and
417 Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. *Communications*
418 *of the ACM*, 65(1):99–106, 2021. 6
- 419 Seungjun Nah, Tae Hyun Kim, and Kyoung Mu Lee. Deep multi-scale convolutional neural network
420 for dynamic scene deblurring. In *CVPR*, July 2017. 4
- 421 Eyal Naor, Itai Antebi, Shai Bagon, and Michal Irani. Combining internal and external constraints
422 for unrolling shutter in videos. In *Computer Vision–ECCV 2022: 17th European Conference, Tel*
423 *Aviv, Israel, October 23–27, 2022, Proceedings, Part XVII*, pp. 119–134. Springer, 2022. 1
- 424 Liyuan Pan, Cedric Scheerlinck, Xin Yu, Richard Hartley, Miaomiao Liu, and Yuchao Dai. Bringing
425 a blurry frame alive at high frame-rate with an event camera. In *Proceedings of the IEEE/CVF*
426 *Conference on Computer Vision and Pattern Recognition*, pp. 6820–6829, 2019. 3
- 427 Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito,
428 Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in
429 pytorch. 2017. 6
- 430 Cedric Scheerlinck, Henri Rebecq, Timo Stoffregen, Nick Barnes, Robert Mahony, and Davide
431 Scaramuzza. Ced: Color event camera dataset. In *Proceedings of the IEEE/CVF Conference on*
432 *Computer Vision and Pattern Recognition Workshops*, pp. 0–0, 2019. 6, 14
- 433 Wei Shang, Dongwei Ren, Dongqing Zou, Jimmy S Ren, Ping Luo, and Wangmeng Zuo. Bringing
434 events into video deblurring with non-consecutively blurry frames. In *Proceedings of the IEEE/CVF*
435 *International Conference on Computer Vision*, pp. 4531–4540, 2021. 3
- 436 Vincent Sitzmann, Julien Martel, Alexander Bergman, David Lindell, and Gordon Wetzstein. Im-
437 plicit neural representations with periodic activation functions. *Advances in Neural Information*
438 *Processing Systems*, 33:7462–7473, 2020. 3
- 439 Chen Song, Qixing Huang, and Chandrajit Bajaj. E-cir: Event-enhanced continuous intensity recovery.
440 In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp.
441 7803–7812, 2022. 2, 3
- 442 Chen Song, Chandrajit Bajaj, and Qixing Huang. Deblursr: Event-based motion deblurring under the
443 spiking representation. *arXiv preprint arXiv:2303.08977*, 2023. 2, 3, 7, 8, 20
- 444 Shuochen Su and Wolfgang Heidrich. Rolling shutter motion deblurring. In *Proceedings of the IEEE*
445 *Conference on Computer Vision and Pattern Recognition*, pp. 1529–1537, 2015. 1
- 446 Deqing Sun, Xiaodong Yang, Ming-Yu Liu, and Jan Kautz. Pwc-net: Cnns for optical flow using
447 pyramid, warping, and cost volume. In *Proceedings of the IEEE conference on computer vision*
448 *and pattern recognition*, pp. 8934–8943, 2018. 5
- 449 Lei Sun, Christos Sakaridis, Jingyun Liang, Qi Jiang, Kailun Yang, Peng Sun, Yaozu Ye, Kaiwei
450 Wang, and Luc Van Gool. Event-based fusion for motion deblurring with cross-modal attention. In
451 *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022,*
452 *Proceedings, Part XVIII*, pp. 412–428. Springer, 2022. 3

- 453 Stepan Tulyakov, Daniel Gehrig, Stamatios Georgoulis, Julius Erbach, Mathias Gehrig, Yuanyou Li,
454 and Davide Scaramuzza. Time lens: Event-based video frame interpolation. In *Proceedings of the*
455 *IEEE/CVF conference on computer vision and pattern recognition*, pp. 16155–16164, 2021. [2](#), [7](#),
456 [8](#), [9](#), [15](#), [17](#), [18](#)
- 457 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz
458 Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing*
459 *systems*, 30, 2017. [6](#), [9](#)
- 460 Bishan Wang, Jingwei He, Lei Yu, Gui-Song Xia, and Wen Yang. Event enhanced high-quality image
461 recovery. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August*
462 *23–28, 2020, Proceedings, Part XIII 16*, pp. 155–171. Springer, 2020. [2](#), [3](#), [5](#), [7](#), [8](#), [13](#)
- 463 Wenhai Wang, Jifeng Dai, Zhe Chen, Zhenhang Huang, Zhiqi Li, Xizhou Zhu, Xiaowei Hu, Tong Lu,
464 Lewei Lu, Hongsheng Li, et al. Internimage: Exploring large-scale vision foundation models with
465 deformable convolutions. *arXiv preprint arXiv:2211.05778*, 2022a. [5](#), [13](#), [20](#)
- 466 Zhixiang Wang, Xiang Ji, Jia-Bin Huang, Shin’ichi Satoh, Xiao Zhou, and Yinqiang Zheng. Neural
467 global shutter: Learn to restore video from a rolling shutter camera with global reset feature.
468 In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp.
469 17794–17803, 2022b. [3](#)
- 470 Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from
471 error visibility to structural similarity. *IEEE transactions on image processing*, 13(4):600–612,
472 2004. [6](#)
- 473 Zirui Wang, Shangzhe Wu, Weidi Xie, Min Chen, and Victor Adrian Prisacariu. Nerf-: Neural
474 radiance fields without known camera parameters. *arXiv preprint arXiv:2102.07064*, 2021. [3](#), [6](#)
- 475 Fang Xu, Lei Yu, Bishan Wang, Wen Yang, Gui-Song Xia, Xu Jia, Zhendong Qiao, and Jianzhuang
476 Liu. Motion deblurring with real events. In *Proceedings of the IEEE/CVF International Conference*
477 *on Computer Vision*, pp. 2583–2592, 2021. [3](#)
- 478 Kaihao Zhang, Wenhan Luo, Yiran Zhong, Lin Ma, Bjorn Stenger, Wei Liu, and Hongdong Li.
479 Deblurring by realistic blurring. In *Proceedings of the IEEE/CVF Conference on Computer Vision*
480 *and Pattern Recognition*, pp. 2737–2746, 2020. [4](#)
- 481 Xiang Zhang and Lei Yu. Unifying motion deblurring and frame interpolation with events. In
482 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp.
483 17765–17774, 2022. [3](#), [6](#)
- 484 Xinyu Zhang, Hefei Huang, Xu Jia, Dong Wang, and Huchuan Lu. Neural image re-exposure. *arXiv*
485 *preprint arXiv:2305.13593*, 2023. [3](#)
- 486 Xu Zheng, Yexin Liu, Yunfan Lu, Tongyan Hua, Tianbo Pan, Weiming Zhang, Dacheng Tao, and Lin
487 Wang. Deep learning for event-based vision: A comprehensive survey and benchmarks. *arXiv*
488 *preprint arXiv:2302.08890*, 2023. [13](#)
- 489 Zhihang Zhong, Yinqiang Zheng, and Imari Sato. Towards rolling shutter correction and deblurring
490 in dynamic scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*
491 *Recognition*, pp. 9219–9228, 2021. [1](#), [3](#), [7](#), [8](#)
- 492 Xinyu Zhou, Peiqi Duan, Yi Ma, and Boxin Shi. Evunroll: Neuromorphic events based rolling shutter
493 image correction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*
494 *Recognition*, pp. 17775–17784, 2022. [1](#), [2](#), [3](#), [6](#), [7](#), [8](#), [9](#), [14](#), [15](#), [16](#), [17](#), [18](#)
- 495 Alex Zihao Zhu, Liangzhe Yuan, Kenneth Chaney, and Kostas Daniilidis. Unsupervised event-based
496 learning of optical flow, depth, and egomotion. In *Proceedings of the IEEE/CVF Conference on*
497 *Computer Vision and Pattern Recognition*, pp. 989–997, 2019. [5](#)