# EV-GS: EVENT-BASED GAUSSIAN SPLATTING FOR EFFICIENT AND ACCURATE RADIANCE FIELD RENDERING

*Jingqian Wu, Shuo Zhu, Chutian Wang, Edmund Y. Lam\**

The University of Hong Kong, Pokfulam, Hong Kong SAR, China

## ABSTRACT

Computational neuromorphic imaging (CNI) with event cameras offers advantages such as minimal motion blur and enhanced dynamic range, compared to conventional frame-based methods. Existing event-based radiance field rendering methods are built on neural radiance field, which is computationally heavy and slow in reconstruction speed. Motivated by the two aspects, we introduce Ev-GS, the first CNI-informed scheme to infer 3D Gaussian splatting from a monocular event camera, enabling efficient novel view synthesis. Leveraging 3D Gaussians with pure event-based supervision, Ev-GS overcomes challenges such as the detection of fast-moving objects and insufficient lighting. Experimental results show that Ev-GS outperforms the method that takes frame-based signals as input by rendering realistic views with reduced blurring and improved visual quality. Moreover, it demonstrates competitive reconstruction quality and reduced computing occupancy compared to existing methods, which paves the way to a highly efficient CNI approach for signal processing.

*Index Terms*— Event Camera, Radiance Field Rendering, 3D Gaussian splatting, Computational Neuromorphic Imaging

## 1. INTRODUCTION

The radiance field rendering aims to output the representation of the distribution and intensity of light across 3D space [1]. With advancements in deep learning frameworks and CUDA-based rasterization, this rendering task has become vital for generating dense, photorealistic renderings of scenes in a 3D-consistent manner, playing a pivotal role in both computer vision and computer graphics [2]. Efforts have been directed towards curating real-scene datasets, often reconstructed from photographs or scans, to serve as benchmarks for evaluating algorithm performance on real-world data [3]. However, real-world images captured by frame-based cameras frequently suffer from motion blur, caused by rapid motion and prolonged exposure duration [1, 4]. This phenomenon occurs when each pixel of the camera integrates light from different points in the scene during exposure, resulting in color value blending. Motion blur leads to a loss of information, hindering tasks such as radiance field rendering and subsequent processing [5]. Simply reducing exposure time to mitigate motion blur is often impractical, as it compromises light reception and exacerbates noise levels [6]. Additionally, conventional frame-based cameras exhibit limited dynamic range, causing bright areas to saturate with white and dark areas to lose detail, particularly crucial information like text [7].

In contrast, event cameras capture asynchronous brightness changes per pixel, offering advanced properties such as high temporal resolution, high dynamic range, lower power consumption, and reduced latency compared to conventional frame-based cameras [8]. Computational neuromorphic imaging (CNI) is a novel paradigm to harness the advanced properties of event cameras for numerous applications[9]. In this context, to utilize the advantage of the event camera, we proposed the first CNI-informed method for inferring 3D Gaussian splatting (GS) from a monocular event camera, enabling efficient and accurate novel view synthesis for objects in grayscale space. Specifically, Our approach, named Ev-GS, leverages the advantages by employing 3D Gaussians as a flexible and efficient representation with purely event-based signal supervision, allowing accurate representation of challenging scenes for traditional frame-based supervision (e.g., motion blur, lacked frame number under high-speed movements, or insufficient lighting), while achieving high-quality rendering through faster training and real-time performance, especially for complex scenes and high-resolution outputs. Experiments show that our Ev-GS approach surpasses frame-based cameras by rendering realistic views with reduced blurring and improved visual quality on real-world datasets. Moreover, compared to existing rendering methods on synthetic datasets, Ev-GS demonstrates competitive rendering quality and significant efficiency improvements, including real-time rendering speed, memory occupancy, and training cost. We summarize our technical contribution as follows:

- We introduce EV-GS, marking the first CNI-informed attempt to infer 3D GS from a monocular event camera, which enables a much more efficient synthesis of novel views for objects within the grayscale space compared to other state-of-the-art rendering methods.

* Corresponding author: elam@eee.hku.hk

- We propose a novel event stream utilization and supervision framework specifically designed for differentiable 3D Gaussian-like methods for accurate and realistic scene rendering.

## 2. RELATED WORK

### 2.1. Neural Rendering and Radiance Field

Recently, 3D GS [2] has emerged as a compelling alternative to Neural Radiance Field (NeRF) [10] for 3D representation, exhibiting notable quality and speed improvements across both 3D and 4D reconstruction tasks. Its efficient differentiable rendering implementation and model design streamline training processes without necessitating spatial pruning [2]. Despite its various advantages, a significant challenge lies in acquiring comprehensive and accurate scene information efficiently for training a 3D Gaussian model. Many approaches [2, 11, 12] adopt scene collection methods by capturing videos from a moving camera. This method offers the advantage of efficiently capturing training data: instead of capturing photos from every angle or direction, which can be time-consuming and prone to overlooking information, a single camera mounted on a moving device traverses the entire scene, ensuring all necessary information is captured. However, using frame-based cameras for this purpose may introduce issues such as motion blur and sparse frames, leading to deficiencies in view rendering results [1, 4]. To address this challenge, we present a new method to learn 3D Gaussian scene representations from event streams. It enables dense photorealistic novel view synthesis, overcoming the limitations associated with frame-based cameras.

### 2.2. Neuromorphic Radiance Field Rendering

The distinctive features of event cameras, including their capacity to prevent motion blur, provide a high dynamic range, ensure low latency, and consume minimal power, have spurred their adoption in the realm of computer vision and computational imaging with vital applications [13, 14, 15, 9]. Several approaches have been proposed to address the view synthesis challenge using NeRF [10] with event data [1, 16], leveraging volumetric rendering with either pure event or semi-event (blurred RGB involved) supervision. However, a significant drawback is the time-consuming optimization of NeRF. The computational demands of training and optimizing an event NeRF pipeline, in terms of both training time and GPU memory, are roughly 83 times more than the 3D GS pipeline. Moreover, the utilization of high-dimensional multilayer perception networks in the NeRF architecture results in a slower view-rendering speed of around 190 times, which may pose limitations for real-time rendering applications.

In this study, to the best of our knowledge, we proposed an adaptation of 3D GS [2] to address view synthesis challenges,

marking the first instance of such an application in this context. Our approach achieved efficient training and rendering processes while upholding realistic visual quality.

## 3. METHOD

### 3.1. Preliminary on 3D Gaussian Splatting

3D GS [2] portrays a detailed 3D scene by utilizing point clouds, with Gaussians utilized to delineate the scene's structure. In this depiction, each Gaussian is characterized by a central point, denoted as $\mathbf{x}$, and a covariance matrix $\Sigma$. The central point $\mathbf{x}$ is commonly referred to as the mean value of the Gaussian

$$G(x) = \exp\left(-\frac{1}{2}\mathbf{x}^T \Sigma^{-1}\mathbf{x}\right). \tag{1}$$

For the purpose of differentiable optimization, the covariance matrix $\Sigma$ can undergo decomposition into a rotation matrix $R$ and a scaling matrix $S$

$$\Sigma = RSS^T R^T. \tag{2}$$

To generate renderings from different perspectives, the technique of splatting, as outlined in [17], is employed to position the Gaussians on the camera planes. This method, initially introduced in [18], entails a viewing transformation denoted by $W$ and the Jacobian $J$ of the affine approximation of the projective transformation. Utilizing these parameters, the covariance matrix $\Sigma'$ in camera coordinates is

$$\Sigma' = JW \Sigma W^T J^T. \tag{3}$$

In summary, each Gaussian point within the model is defined by a collection of attributes: its position, represented by $\mathbf{x} \in \mathbb{R}^3$, its color depicted by spherical harmonics coefficients $\mathbf{c} \in \mathbb{R}^k$ (where $k$ represents the degrees of freedom), its opacity $\alpha \in \mathbb{R}$, a rotation quaternion $\mathbf{q} \in \mathbb{R}^4$, and a scaling factor $\mathbf{s} \in \mathbb{R}^3$. Specifically, for each pixel, the color and opacity of all Gaussians are computed based on the Gaussian representation outlined in Equation 1. The blending process for color $C$ of N-ordered points overlapping a pixel adheres to a precise formula

$$C = \sum_{i \in N} \mathbf{c}_i \alpha_i \prod_{j=1}^{i-1}(1 - \alpha_j), \tag{4}$$

where variables $\mathbf{c}_i$ and $\alpha_i$ denote the color and density of a specific point, respectively. These values are influenced by a Gaussian with a covariance matrix $\Sigma$, which is subsequently adjusted by adjustable per-point opacity and spherical harmonics color coefficients.
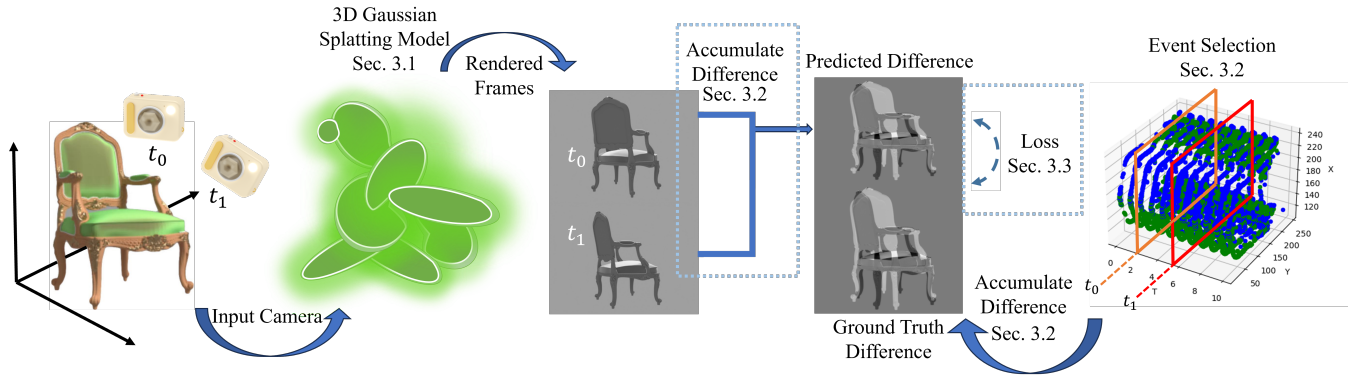
**Fig. 1.** An overview of Ev-GS: a novel method for learning the radiance field volume from a moving event camera via 3D GS (Section 3.1). We establish a link between the observed events and the rendered views at two distinct timestamps, $t_0$ and $t_1$, utilizing an event-based integral (Section 3.2), under the supervision of pure event signal (Section 3.3).

## 3.2. Event Stream Utilization

Each event $e_k$ is described as a tuple $(\mathbf{x}_k, t_k, p_k)$, which occur asynchronously at pixel $\mathbf{x}_k = (x_k, y_k)$ at micro-second timestamp $t_k$.

The polarity $p_k \in \{-1, +1\}$ denotes either an increase or decrease in the logarithmic brightness $L(\mathbf{x}_k, t_k)$ by the contrast threshold $A$. In other words, an event at time $t_k$ is triggered if the following condition is met

$$\Delta L \triangleq L(\mathbf{x}_k, t_k) - L(\mathbf{x}_k, t_{k-1}) = p_k A, \quad (5)$$

where $t_{k-1}$ is the timestamp that the previous event occurs at pixel $u_k$.

As described in Figure 1, our goal is to render radiance field representation from differentiable 3D Gaussian functions under pure event signal supervision, with no RGB or grayscale frame-based data involved. To achieve this goal, we have to formulate the ground truth event data to a differentiable supervision signal and train a 3D GS model to render such representation. The core principle behind our algorithm is to generate two rendering results at two different camera poses on two timestamps, supervised by the ground truth event signal, which is the accumulated event frame between those two timestamps.

Specifically, we set a max window length of $W = 50$, and randomly select a window length of $w \in \text{range}(W)$. For each timestamp $t$, we calculate two associated rendering result from the 3D Gaussian model $I_t = G(c_t)$ and $I_{t-w} = G(c_{t-w})$, where $I_t$ and $I_{t-w}$ are the rendered grayscale frame at time $t$ and $t-w$, $G$ is the 3D GS model, and $c_t$ and $c_{t-w}$ are the camera pose at time $t$ and $t-w$. We represented its logarithmic image as $L(I_t) = \log\left((I_t)^g + \epsilon\right)$, where $\epsilon = 1 \times 10^{-5}$, and $g$ denotes a fixed gamma correction value set to 4.8 across all experiments. It conforms to the grayscale gamma curve, with Gamma 4.8 serving as the recommended smooth approx-

imation [19]. As a result, we derive the predicted accumulative difference $E_{pred} = L(I_t) - L(I_{t-w})$.

On the other hand, to utilize event data and formulate it to a supervision signal for $E_{pred}$, following [20], we aggregate the polarities of all events that transpire between the selected time $t$ and $t-w$ based on their positional information $u$, such that

$$E_{gt} = \int_t^{t+w} e_k(\mathbf{x}_k, p_k, t_k)\, dt, \quad (6)$$

where $E_{gt}$ is the aggregated result.

## 3.3. Event Stream Based Supervision

To effectively supervise $E_{pred}$ from $E_{gt}$, following v2e [21] and [1], we apply the linlog mapping described to derive the predicted logarithmic brightness difference for both $E_{pred}$ and $E_{gt}$. With a normalized $L_2$ loss, the loss $L_e$ can be calculated as

$$L_e(x, y) = \frac{\sum_{i=1}^{H} \sum_{j=1}^{W} (\|L(x)\|_2^2 - \|L(y)\|_2^2)^2}{H \times W}, \quad (7)$$

where for an arbitrary $u$

$$L(u) = \text{linlog}\left(I(u)\right) \triangleq \begin{cases} I(u) \times \ln(B)/B & I(u) < B \\ \ln\left(I(u)\right) & \text{otherwise.} \end{cases} \quad (8)$$

The threshold $B$ delineates the linear region, where no logarithmic mapping is applied. The value of 20 is used for $B$ in all our experiments. $x$ and $y$ are $E_{pred}$, $E_{gt}$ correspondingly.

We also kept the D-SSIM loss used in the original 3D GS article, as we found a small coefficient would improve the

rendering result. Following previous works, we calculate the D-SSIM loss as

$$L_{SSIM}(x,y) = \frac{(2\mu_x\mu_y + c_1)(2\sigma_{xy} + c_2)}{(\mu_x^2 + \mu_y^2 + c_1)(\sigma_x^2 + \sigma_y^2 + c_2)}, \quad (9)$$

where $\mu$ is the pixel sample mean, $\sigma_x$ is the variance of x, $\sigma_{xy}$ is the covariance of $x$ and $y$, $c_1 = (k_1 L)^2$, $c_2 = (k_2 L)^2$ are two variables to stabilize the division with weak denominator, L is the window range of the pixel values, and $k_1$, $k_2$ are constants. $x$ and $y$ are $E_{pred}, E_{gt}$ correspondingly. Therefore, the final loss function can be derived as follows

$$L_{total} = L_e + \lambda(1 - L_{SSIM}), \quad (10)$$

where $\lambda$ is 0.1 for all our experiments.

## 4. EXPERIMENTS

### 4.1. Overview

We analyze a total of four synthetic sequences and four real sequences. For the synthetic sequences, we utilize the 3D models sourced from [10]. Each scene undergoes rendering for a duration of one second, depicting a 360-degree rotation of the camera around the object at $1000$ fps (frames per second), resulting in 1000 RGB images, then turned into 1000 grayscale images. From these images, we generate the event stream using the model [22], with the corresponding camera intrinsics and extrinsic directly applied in our approach.

In our real-data experiments, we capture footage of four objects using the DAVIS 346C event camera. It is worth noting the challenges in creating a stable and calibrated setup in real-world scenarios where the event camera rotates around an object. To address this, we maintain the camera's static position and place the objects on a turntable with a maximum turning speed of five seconds per round. In this setup, maintaining constant lighting irrespective of the object's rotation angle is crucial. To achieve this, we mount a single USB ring light above the object. A photograph of the setup is provided by Figure 2.

The evaluation metrics for the rendering results of 3D GS include the peak signal-to-noise ratio (PSNR) and the structural similarity index (SSIM). PSNR measures signal fidelity by comparing the maximum signal power to noise power, while SSIM assesses image similarity based on luminance, contrast, and structure. These metrics provide quantitative insights into the quality and fidelity of the rendered images.

We further conduct an ablation experiment that demonstrates the necessity of our method design.

### 4.2. Implementing Details

Our implementation is constructed upon the 3D GS codebase [2], leveraging its foundational structure and functionalities. Because 3D GS requires a point cloud as input, we randomly
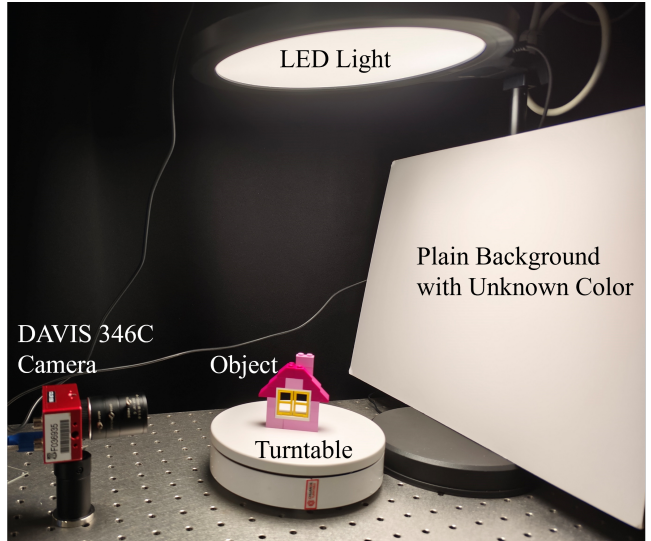


Fig. 2. Demonstration of Hardware Setup.

initialized $10^6$ points to form the initial point cloud since the structure-from-motion [23] initialization, which is originally used in the model, is not applicable for event data. All experiments detailed in this paper are conducted by utilizing the computational resources of an NVIDIA RTX 3090 GPU.

We follow the settings of several hyperparameters for performance and optimization. The total number of iterations during the training process is set to $50,000$. For position optimization, the initial and final learning rates are $1.6 \times 10^{-4}$ and $1.6 \times 10^{-6}$. Learning rates for feature, opacity, scaling, and rotation optimization are $2.5 \times 10^{-3}, 5 \times 10^{-2}, 5 \times 10^{-3}$, and $1 \times 10^{-3}$, respectively.

### 4.3. Synthetic Sequences

As previously mentioned, we evaluate synthetic sequences sourced from Mildenhall et al. [10], encompassing various effects including chairs, hot dogs, ficus, and microphones. Table 1 demonstrates the quantitative comparison EventNeRF [16]. According to our experiment results in Table 1, our method demonstrates significant advantages over EventNeRF across various scenes. In most scenes, our approach achieves a higher PSNR and SSIM than EventNeRF, indicating superior reconstruction fidelity and better structural similarity between the reconstructed and ground truth images. Notably, our method boasts significantly reduced training times, taking only around 10 minutes compared to EventNeRF's 14 hours for all scenes. This accelerated training time allows for more efficient model development and experimentation. Moreover, our method achieves higher temporal resolution, reaching around $60$ fps compared to EventNeRF's $0.32$ fps. A high frame rate is crucial for smooth and responsive visual experiences as it ensures real-time responsiveness, reduces motion sickness, and enhances productivity and com-
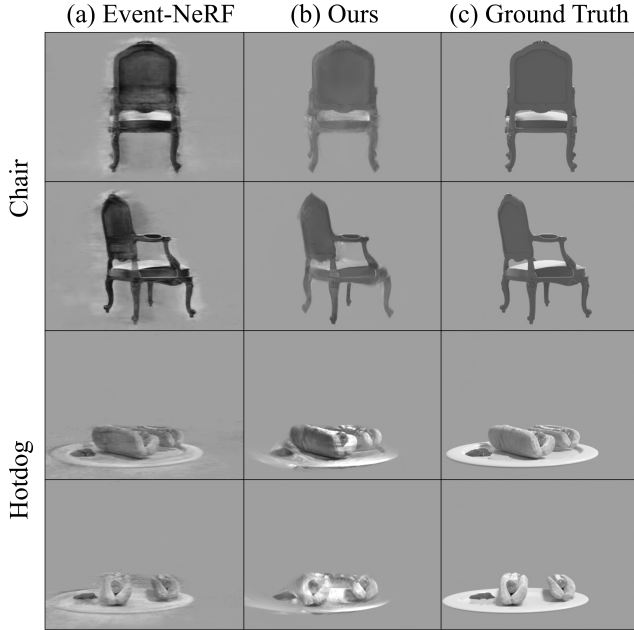
(a) Event-NeRF    (b) Ours    (c) Ground Truth

**Fig. 3**. Visual Comparison between our approach and Event-NeRF.

**Table 1**. Quantitative Comparison Against Event-NeRF [16].

| Metirc | PSNR↑ | SSIM↑ | Training Time↓ | FPS↑ | Memory↓ |
|--------|-------|-------|----------------|------|---------|
| Scene | | | Chair | | |
| EventNeRF | 25.6 | 0.91 | 14h | 0.32 | 15GB |
| Ours | **28.1** | **0.93** | **9min** | **53.1** | **5GB** |
| Scene | | | Ficus | | |
| EventNeRF | 27.1 | 0.91 | 14h | 0.33 | 15GB |
| Ours | **28.1** | **0.92** | **9min** | **63.0** | **5GB** |
| Scene | | | Hotdog | | |
| EventNeRF | 26.0 | 0.92 | 14h | 0.32 | 15GB |
| Ours | 25.7 | **0.93** | **12min** | **55.3** | **5GB** |
| Scene | | | Mic | | |
| EventNeRF | 25.0 | 0.915 | 14h | 0.32 | 15GB |
| Ours | 24.5 | **0.92** | **11min** | **66.02** | **5GB** |

petitiveness. Lastly, our method consumes less memory during the rendering stage. Ev-Gs utilizes only 5 GB compared to EventNeRF's 15 GB, which is advantageous for memory-constrained environments and facilitates scalability.

Qualitatively, we also compare visual results for two sequences. As shown in Fig 3, our Ev-GS effectively learns view-dependent effects, structures, and intensity of the frame compared to the previous method.

### 4.4. Real Sequences

As described in Section 4.1, for real-world sequence, we analyzed a Lego toy scene captured on a turntable. Since no ground-truth RGB data is available, we assess the results visually, as depicted in Figure 4. In scenarios with fast-moving
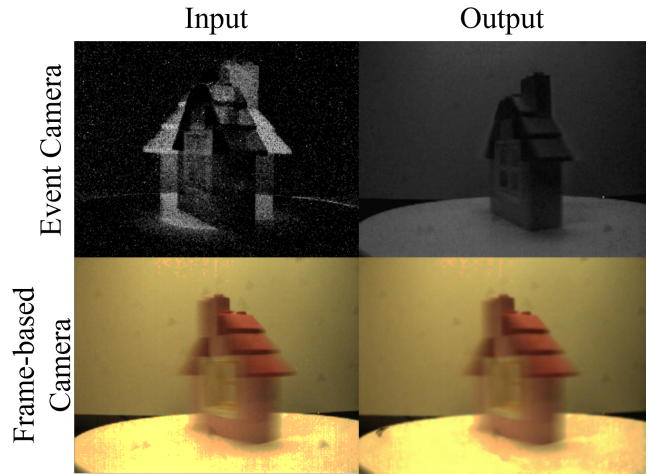


Input    Output

**Fig. 4**. Visual Comparison of real-world data between event signal supervision and frame-based supervision.

cameras, traditional frame-based systems suffer from motion blur, resulting in blurry renderings. Our method, solely reliant on event data, circumvents motion blur and low-contrast issues, ensuring more precise and higher-quality rendering outcomes.

## 5. CONCLUSION

In this work, we introduced Ev-GS, a novel method for inferring 3D Gaussian splatting from monocular event cameras. Leveraging the unique advantages of the CNI paradigm, Ev-GS enables efficient and accurate novel view synthesis in grayscale space. Our approach overcomes challenges by employing purely event-based supervision, resulting in superior rendering quality compared to existing methods. Experimental results demonstrate the effectiveness of Ev-GS in rendering realistic views with reduced blurring and improved visual quality on real-world datasets. Moreover, our method shows significant efficiency improvements, including real-time reconstruction speed and reduced memory occupancy, highlighting its potential for various applications within the CNI framework. However, Ev-GS still suffers from reconstructing difficult scenes, especially in complex objects with challenging textures. Future work will be done to fill the gap to enhance the robustness of reconstruction results both quantitatively and qualitatively. Overall, we believe Ev-GS will provide an enlightening reference and shed light on signal processing with event representation and its optical applications.

## 6. REFERENCES

[1] Simon Klenk, Lukas Koestler, Davide Scaramuzza, and Daniel Cremers, "E-NeRF: Neural radiance fields from a moving event camera," *IEEE Robotics and Automation Letters*, vol. 8, no. 3, pp. 1587–1594, 2023.

[2] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis, "3D Gaussian splatting for real-time radiance field rendering," *ACM Transactions on Graphics (TOG)*, vol. 42, no. 4, pp. 1–14, 2023.

[3] Ziyi Yang, Xinyu Gao, Wen Zhou, Shaohui Jiao, Yuqing Zhang, and Xiaogang Jin, "Deformable 3D Gaussians for high-fidelity monocular dynamic scene reconstruction," *arXiv preprint arXiv:2309.13101*, 2023.

[4] Javier Hidalgo-Carrió, Guillermo Gallego, and Davide Scaramuzza, "Event-aided direct sparse odometry," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 5781–5790.

[5] Richard Szeliski, *Computer vision: algorithms and applications*, 2022.

[6] Paul Debevec, Tim Hawkins, Chris Tchou, Haarm-Pieter Duiker, Westley Sarokin, and Mark Sagar, "Acquiring the reflectance field of a human face," in *Conference on Computer Graphics and Interactive Techniques*, 2000, pp. 145–156.

[7] Boyoon Jung and Gaurav S Sukhatme, "Real-time motion tracking from a mobile robot," *International Journal of Social Robotics*, vol. 2, pp. 63–78, 2010.

[8] Yi-Fan Zuo, Jiaqi Yang, Jiaben Chen, Xia Wang, Yifu Wang, and Laurent Kneip, "Devo: Depth-event camera visual odometry in challenging conditions," in *2022 International Conference on Robotics and Automation (ICRA)*, 2022, pp. 2179–2185.

[9] Shuo Zhu, Chutian Wang, Haosen Liu, Pei Zhang, and Edmund Y Lam, "Computational neuromorphic imaging: Principles and applications," in *Computational Optical Imaging and Artificial Intelligence in Biomedical Sciences*, 2024, vol. 12857, SPIE, pp. 4–10.

[10] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng, "NeRF: Representing scenes as neural radiance fields for view synthesis," *Communications of the ACM*, vol. 65, no. 1, pp. 99–106, 2021.

[11] Yunzhi Yan, Haotong Lin, Chenxu Zhou, Weijie Wang, Haiyang Sun, Kun Zhan, Xianpeng Lang, Xiaowei Zhou, and Sida Peng, "Street Gaussians for modeling dynamic urban scenes," *arXiv preprint arXiv:2401.01339*, 2024.

[12] Zhe Li, Zerong Zheng, Lizhen Wang, and Yebin Liu, "Animatable Gaussians: Learning pose-dependent Gaussian maps for high-fidelity human avatar modeling," *arXiv preprint arXiv:2311.16096*, 2023.

[13] Pei Zhang, Haosen Liu, Zhou Ge, Chutian Wang, and Edmund Y. Lam, "Neuromorphic imaging with joint image deblurring and event denoising," *IEEE Transactions on Image Processing*, vol. 33, pp. 2318–2333, March 2024.

[14] Shuo Zhu, Zhou Ge, Chutian Wang, Jing Han, and Edmund Y Lam, "Efficient non-line-of-sight tracking with computational neuromorphic imaging," *Optics Letters*, vol. 49, no. 13, pp. 3584–3587, 2024.

[15] Chutian Wang, Shuo Zhu, Pei Zhang, Jianqing Huang, Kaiqiang Wang, and Edmund Y. Lam, "Neuromorphic shack-hartmann wave normal sensing," *arXiv preprint arXiv:2404.15619*, 2024.

[16] Viktor Rudnev, Mohamed Elgharib, Christian Theobalt, and Vladislav Golyanik, "EventNeRF: Neural radiance fields from a single colour event camera," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 4992–5002.

[17] Wang Yifan, Felice Serena, Shihao Wu, Cengiz Öztireli, and Olga Sorkine-Hornung, "Differentiable surface splatting for point-based geometry processing," *ACM Transactions on Graphics (TOG)*, vol. 38, no. 6, pp. 1–14, 2019.

[18] Matthias Zwicker, Hanspeter Pfister, Jeroen Van Baar, and Markus Gross, "Surface splatting," in *Conference on Computer Graphics and Interactive Techniques*, 2001, pp. 371–378.

[19] International Electrotechnical Commission et al., "Multimedia systems and equipment-color measurement and management-part 2-1," *Color management-Default RGB color space-sRGB*, 1999.

[20] Alex Zihao Zhu, Liangzhe Yuan, Kenneth Chaney, and Kostas Daniilidis, "Unsupervised event-based optical flow using motion compensation," in *European Conference on Computer Vision (ECCV) Workshops*, 2018, pp. 0–0.

[21] Yuhuang Hu, Shih-Chii Liu, and Tobi Delbruck, "v2e: From video frames to realistic dvs events," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 1312–1321.

[22] Daniel Gehrig, Mathias Gehrig, Javier Hidalgo-Carrió, and Davide Scaramuzza, "Video to events: Recycling video datasets for event cameras," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 3586–3595.

[23] Johannes L Schonberger and Jan-Michael Frahm, "Structure-from-motion revisited," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 4104–4113.