

# Keypoint-based Dynamic Object 6-DoF Pose Tracking via Event Camera

Zhe Wang, Qijin Song, Zihao Li, Jingyu Xiao, and Weibang Bai\*

**Abstract**—Accurate 6-DoF pose estimation of objects is critical for robots to perform precise manipulation tasks. However, for dynamic object pose estimation, conventional camera-based approaches face several major challenges, such as motion blur, sensor noise, and low-light limitation. To address these issues, we employ event cameras, whose high dynamic range and low latency offer a promising solution. Furthermore, we propose a keypoint-based detection and tracking approach for dynamic object pose estimation. Firstly, a keypoint detection network is constructed to extract keypoints from the time surface generated by the event stream. Subsequently, the polarity and spatial coordinates of the events are leveraged, and the event density in the vicinity of each keypoint is utilized to achieve continuous keypoint tracking. Finally, a hash mapping is established between the 2D keypoints and the 3D model keypoints, and the EPnP algorithm is employed to estimate the 6-DoF pose. Experimental results demonstrate that, whether in simulated or real event environments, the proposed method outperforms the event-based state-of-the-art methods in terms of both accuracy and robustness.

## I. INTRODUCTION

Object pose estimation aims to compute the precise position and orientation of a target in the world coordinate system, thereby obtaining its complete 6-DoF representation in three-dimensional space. In the field of robotics, object pose estimation plays a foundational role. It provides prior pose information for autonomous grasping and placement, enabling robotic grippers to align precisely with target objects [1]. For tool operation and interaction, pose estimation helps robots infer tool geometry and orientation to perform actions such as cutting, rotating, and plugging [2]. In automated assembly and manufacturing, accurate pose estimation ensures precise part alignment on high-speed production lines, thereby enhancing stability and efficiency [3].

In recent years, object pose estimation has made significant progress, and the types of visual sensors employed have become increasingly diverse, including monocular cameras [4], [5], [6], stereo cameras [7], [8], [9], and depth cameras [10], [11]. Do et al. [5] introduced LieNet, a real-time framework that detects, segments, and estimates 6D object poses

This work is supported by the Shanghai Pujiang Program under grant 23PJ1408500, and the MoE Key Laboratory of Intelligent Perception and Human-Machine Collaboration (KLIP-HuMaCo). The experiments of this work were supported by the Core Facility Platform of Computer Science and Communication, ShanghaiTech University. (\*Corresponding author: Weibang Bai (Email: wbbai@shanghaitech.edu.cn).

Zhe Wang, Qijin Song, Zihao Li and Weibang Bai are with the ShanghaiTech Automation and Robotics (STAR) Center, School of Information Science and Technology, ShanghaiTech University, Shanghai, 201210, China  
Jingyu Xiao is with WLSA Shanghai Academy, and is an intern with the School of Information Science and Technology, ShanghaiTech University, Shanghai, 201210, China

from a single RGB image using a rotation representation based on the Lie algebra. Franke et al. [8] proposed 6D-Vision, a Kalman filter-based stereo-motion fusion for real-time 3D position and motion estimation of image points, enabling robust obstacle detection. These sensors have enabled a wide range of object 6-DoF pose estimation applications under static or slightly dynamic conditions [5], [8]. However, the limited spatiotemporal resolution of these cameras degrades performance, making them unsuitable for accurate 6-DoF pose estimation of highly dynamic objects [12].

The event camera [13], [14] is a kind of bio-inspired sensor producing asynchronous events when the illumination of a single pixel changes, which makes it particularly effective in capturing detailed information of objects in high-speed motion. Therefore, it is well-suited for 6-DoF pose estimation of highly dynamic objects. In recent years, numerous event-based 6-DoF pose algorithms have been proposed.

However, existing event-based object 6-DoF pose tracking methods still exhibit significant limitations: on one hand, they show insufficient adaptability to non-planar geometries [15], [16]; on the other hand, they typically rely on a predefined and aligned initial pose to enable subsequent tracking [15], [17]. However, in service and industrial assembly environments, most components to be handled are curved-surface objects, for which existing methods show limited adaptability. Moreover, it is often impractical to predefine and align the initial pose of the target object in real-world scenarios, thereby imposing constraints on subsequent pose tracking.

To address these challenges, we introduce a keypoint-based pipeline for object detection and pose tracking with event cameras. The method employs keypoint detection for tracking 6-DoF pose of curved-surface objects and automatically initializes object pose via 2D-3D correspondences. The overall pipeline is as follows: Firstly, we introduce a lightweight network for efficient multi-scale feature extraction and robust event-based keypoints detection. Subsequently, to achieve 6-DoF pose tracking under highly dynamic conditions, this work proposes an event-based density keypoints tracking algorithm, in which temporal stability is enhanced by introducing an extended Kalman filter that effectively reduces drift and jitter. In addition, we adopt a structure-sensitive loss that jointly enforces heatmap, coordinate, and geometric constraints, thereby ensuring accurate keypoint localization under sparse event conditions.

The main contributions of this work are as follows:

- We present a keypoint-based event-camera pipeline for 6-DoF pose tracking of dynamic objects, achieving

stable and accurate estimation under high-speed motion.

- We propose a lightweight event-based keypoint detector that reduces target loss from delayed initial 6-DoF estimation while maintaining precise localization under sparse events.
- We propose an event-based density keypoint tracking method with polarity-adaptive matching and EKF for robust tracking under high-speed motion.

## II. RELATED WORK

Object 6-DoF pose estimation and tracking using different vision sensors has been a major focus of research over the past decades. We briefly revisit the two primary categories, including those that utilize traditional RGB cameras and those that employ event cameras.

**Object Pose Tracking through Traditional RGB Cameras:** Recent research on RGB camera-based pose tracking applies edge extraction, feature selection, and image region processing, along with direct optimization and deep learning strategies. Object pose estimation methods based on feature detection achieve accurate spatial pose estimation by extracting line features and keypoint features of the object. He et al. [18] propose ContourPose, a monocular 6-DoF pose estimation method that detects keypoints on object contours to establish 2D–3D correspondences for robust pose prediction of reflective, textureless metal parts.

In addition, deep learning methods have shown impressive performance in object pose estimation. Chen et al. [19] propose EPro-PnP, a probabilistic and differentiable PnP solver that enables end-to-end learning of object pose from monocular images via 2D–3D correspondences. However, such data-driven methods struggle to achieve real-time performance, making it difficult to estimate the pose of fast-moving objects.

Due to hardware limitations, RGB cameras are prone to significant motion blur when capturing fast-moving objects, which limits their suitability for pose estimation in highly dynamic scenarios.

**Object Pose Tracking through Event Cameras:** With their high temporal resolution and low latency, event cameras offer a distinct advantage in tracking 6-DoF pose of fast-moving objects. The first event-camera-based algorithm for 6-DoF object pose estimation and tracking is proposed by D Reverter Valeiras [20]. In this work, they give object model and object’s init pose information, and the algorithm estimates and tracks moving objects by associating incoming events with the 3D structure of the objects. Subsequently, Valeiras et al. [21] present an event-based PnP algorithm to estimate object poses and enable continuous tracking, formulated as a least-squares optimization problem. Liu et al. [16] propose a stereo event-based 6-DoF pose tracking method for uncooperative spacecraft, leveraging geometric associations between asynchronous events and 3D object models. Without relying on deep learning, their approach employs a sparse PnP algorithm to achieve robust and low-latency pose estimation suitable for resource-constrained space applications.

Recent research in machine learning has sparked interest in utilizing data-driven algorithms for event-based object pose tracking. Rathinam et al. [22] introduced SPADES, a realistic and diverse spacecraft pose estimation dataset captured using event cameras under high dynamic lighting conditions. They designed a three-channel event image encoding and occlusion mask filtering strategy, providing a high-quality benchmark for training and evaluating event-driven deep learning models while bridging the sim-to-real domain gap. Yishi et al. [23] proposed a 6-DoF pose estimation algorithm for non-cooperative targets by fusing monocular images with event streams. By constructing event images and employing a cross-modal Transformer architecture, their approach leverages the high dynamic range and motion robustness of event data together with the texture richness of RGB images, achieving high accuracy in challenging environments.

## III. METHOD

In this article, we propose a new event-based technique for tracking the 6-DoF pose of an object. The 6-DoF pose of the object is expressed as a homogeneous transformation matrix  $P = [R \mid T] \in \mathbb{SE}(3)$ , comprising a rotation component  $R \in \mathbb{SO}(3)$  and a translation component  $T \in \mathbb{R}^3$ .

To track the 6-DoF pose of an object, we employ a four stage pipeline, as illustrated in Fig.1. The first stage uses a neural network to accurately detect the keypoints of the object from the time surface derived from the event stream. The second stage matches the 2D keypoints detected by the event camera with the 3D keypoints of the model using a hash table. The third stage leverages event information to compute the event density and integrates a Kalman filter to track keypoints. The fourth stage involves establishing correspondences between the detected 2D keypoints and the object’s 3D coordinates, followed by pose tracking via the EPnP algorithm.

### A. Problem formulation

Unlike conventional cameras, which capture frames at fixed intervals, event cameras output asynchronous streams of events. These cameras feature pixels that operate independently and respond to changes in the logarithmic photocurrent  $L \doteq \log(I)$ . Events can be mathematically described as  $e_k = (\mathbf{x}_k, t_k, p_k)$ , where  $\mathbf{x}_k = (x_k, y_k)^\top$  denotes the location of the pixel,  $t_k$  is the timestamp of events, and  $p_k \in \{+1, -1\}$  indicates the polarity - with +1 representing a positive event (ON) and -1 indicating a negative event (OFF) [27]. An event is triggered when the brightness change at a single pixel exceeds a predefined threshold.

$$\Delta L(\mathbf{x}_k, t_k) \doteq L(\mathbf{x}_k, t_k) - L(\mathbf{x}_k, t_k - \Delta t_k) \quad (1)$$

An event is generated when the change in logarithmic brightness at a pixel reaches a predefined contrast threshold. Specifically, the brightness change associated with an event can be expressed as

$$\Delta L(\mathbf{x}_k, t_k) = p_k D \quad (2)$$

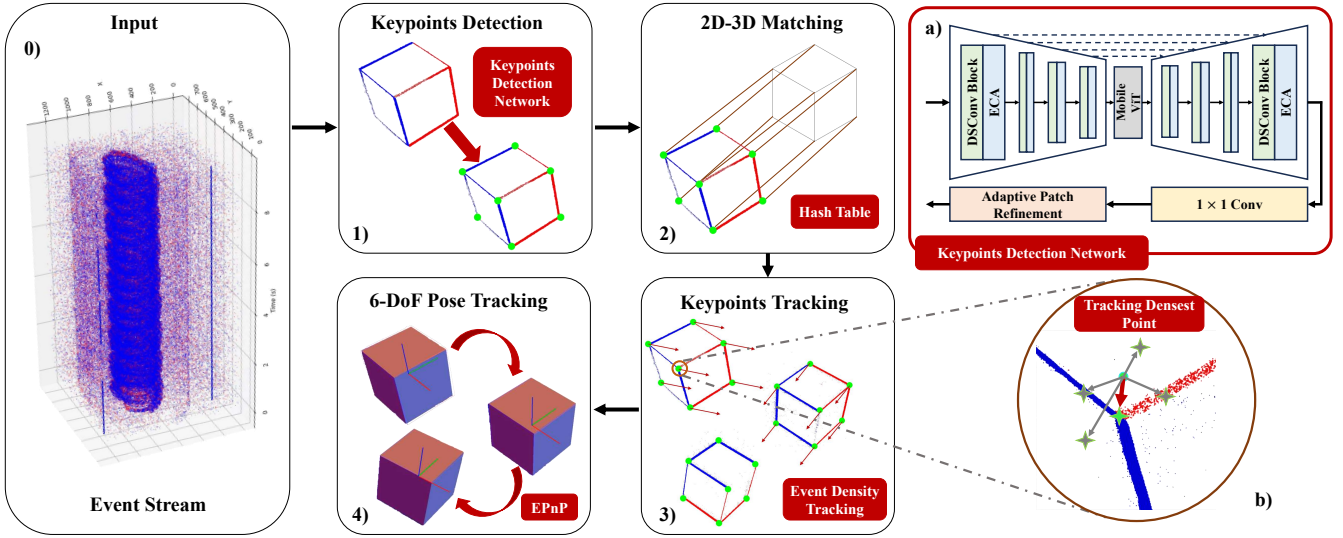


Fig. 1. Overview of the proposed architecture. The proposed framework consists of seven building blocks: 0) obtain event stream of detection object, 1) a keypoints detection network to detect object's keypoints, 2) 2D-3D matching part using hash table, 3) event density tracking part for tracking keypoints by event density information, and 4) object 6-DoF pose tracking by EPnP algorithm. a) is the architecture of this keypoints detection network integrating Depthwise Separable Convolution (DSCConv) [24], Efficient Channel Attention (ECA) [25], MobileViT [26] and Adaptive Patch Refinement module. b) is a local enlarged view of step 3, illustrating the event density tracking process by highlighting the densest point used for keypoint tracking.

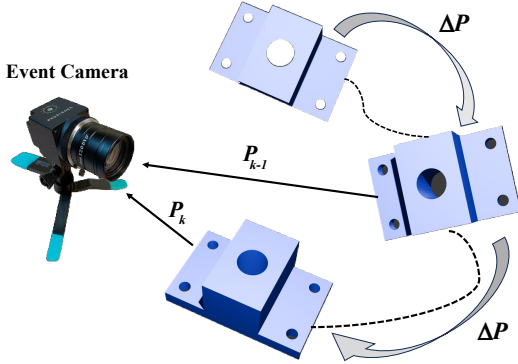


Fig. 2. The geometric interpretation of 6-DoF object pose tracking.  $P_{k-1}$  and  $P_k$  represent the pose of the object at the previous time  $t_{k-1}$  and the current time  $t_k$ , respectively.  $\Delta P$  denotes the change in pose from  $P_{k-1}$  to  $P_k$ .

where  $p_k \in \{+1, -1\}$  denotes the polarity of the event, and  $D$  is the contrast threshold of the event camera. A positive polarity with  $p_k = +1$  corresponds to an increase in brightness, whereas a negative polarity with  $p_k = -1$  corresponds to a decrease in brightness.

In event-based pose tracking, the objective is to continuously estimate the pose  $\mathbf{P}$  of an object from the asynchronous stream of events, thereby determining its position and orientation relative to the event camera over time. Inspired by the concept of recursive pose estimation, the current pose of the object at time  $t_k$  can be inferred based on the previously estimated pose  $\mathbf{P}_{t_{k-1}}$  at time  $t_{k-1}$ , utilizing the information provided by the incoming events at  $t_k$ .

$$\mathbf{P}_{t_k} = \mathbf{P}_{t_{k-1}} \Delta \mathbf{P}^{-1} \quad (3)$$

As illustrated in Fig. 2, the geometric interpretation of 6-DoF object pose tracking indicates that the pose variation between adjacent time steps, denoted as  $\Delta \mathbf{P}$ , is typically

small. Therefore, we initialize the current pose estimation  $\mathbf{P}_{T_k}$  using the previous pose  $\mathbf{P}_{T_{k-1}}$ , and iteratively refine it to obtain an accurate estimation of the current pose.

### B. Event-based keypoints detection

Because event streams provide sparse spatio-temporal information. This makes it difficult to accurately identify the semantic keypoints of object. To address this issue, we design this network to enable effective extraction of semantic keypoint information, as illustrated in Fig. 1-a.

This lightweight network integrates a UNet-based architecture with an Adaptive Patch Refinement module, thereby enabling highly accurate and robust keypoint detection. The proposed keypoints detection network takes the time surface image of the event stream  $I \in \mathbb{R}^{H \times W}$  as input and predicts a set of keypoint heatmaps  $\{H^k\}_{k=1}^K$ , where each heatmap  $H^k \in \mathbb{R}^{H \times W}$  represents the probability response distribution of the  $k$ -th keypoint in the time surface image. The final predicted coordinate  $\hat{\mathbf{p}}^k \in \mathbb{R}^2$  for each keypoint is obtained by extracting the peak response location from the corresponding heatmap:

$$\hat{\mathbf{p}}^k = \arg \max_{(i,j)} H^k(i,j) \quad (4)$$

The overall network architecture is built upon the U-Net framework, consisting of a downsampling encoder, a semantic enhancement bottleneck, and a symmetric upsampling decoder. The encoder employs Depthwise Separable Convolution (DSCConv) [24] combined with Efficient Channel Attention (ECA) [25] to reduce computational cost.

The intermediate layers incorporate the MobileViT [26] module to construct global semantic representations. The intermediate feature map  $F \in \mathbb{R}^{C \times H' \times W'}$  is divided into several non-overlapping patches of size  $r \times r$ . Each patch is flattened into a token sequence  $t_i \in \mathbb{R}^{r^2 \cdot C}$ , which serves

as input to the transformer. After processing by multiple layers of the transformer encoder, the output is reshaped to its spatial dimensions and fused by convolution to reconstruct the original feature  $F' \in \mathbb{R}^{C \times H' \times W'}$ . This module preserves local perceptual ability while establishing long-range dependencies across patches, effectively alleviating structural discontinuities in the event representation.

To further improve the accuracy and sharpness of the heatmap responses, we design an Adaptive Patch Refinement module. Specifically, a local patch is extracted from the coarse heatmap around the  $k$ -th keypoint:

$$\mathcal{P}_k = H^k[x_k - r : x_k + r, y_k - r : y_k + r] \quad (5)$$

Here,  $(x_k, y_k) = \hat{p}_k$  denotes the predicted keypoint location and the patch size is  $r \times r$ . The notation  $[\cdot]$  represents a slicing operation that extracts a local patch from the heatmap centered at  $(x_k, y_k)$  with a radius  $r$ . Each patch is fed into a Tinyu-Net [28] module for local structural modeling and enhancement. The refined patch is then stitched back to the original heatmap location to improve boundary clarity and peak localization stability.

With the above architecture, the network possesses multi-scale perception capabilities from global to local levels. It can effectively capture both the geometric structure and response center of keypoints, which is particularly suitable for high-speed motion and sparse-event scenes in event-based vision. Keypoints detection provides the initial position information for keypoints tracking, while continuous tracking of keypoints constitutes the necessary condition for achieving 6-DoF pose tracking.

### C. Loss function

To improve both the accuracy and geometric consistency of event-based keypoints detection, we propose a Structure-Aware Keypoint Heatmap Loss, which jointly supervises the heatmap response quality, keypoint localization accuracy, and spatial structure stability. The overall loss function is formulated as a weighted combination of heatmap loss, coordinate loss, and structure preservation loss, which can be presented as:

$$\mathcal{L}_{\text{SAKHL}} = \lambda_1 \mathcal{L}_{\text{heatmap}} + \lambda_2 \mathcal{L}_{\text{coord}} + \lambda_3 \mathcal{L}_{\text{structure}} \quad (6)$$

where  $\lambda_1, \lambda_2, \lambda_3$  are the weights controlling the contribution of each term.

The heatmap loss function  $\mathcal{L}_{\text{heatmap}}$  measures the pixel-wise mean squared error between the predicted heatmaps  $\hat{H}_{b,k}$  and the ground-truth heatmaps  $H_{b,k}$ , in each batch size  $b$ :

$$\mathcal{L}_{\text{heatmap}} = \frac{1}{BKH W} \sum_{b=1}^B \sum_{k=1}^K \sum_{i=1}^H \sum_{j=1}^W (\hat{H}_{b,k}(i, j) - H_{b,k}(i, j))^2 \quad (7)$$

The coordinate loss  $\mathcal{L}_{\text{coord}}$  constrains the predicted keypoint locations  $\hat{p}_{b,k}$ , derived from the heatmap peaks, to remain close to the labeled ground-truth coordinates  $p_{b,k}$  by applying the  $\ell_1$  norm:

$$\mathcal{L}_{\text{coord}} = \frac{1}{BK} \sum_{b=1}^B \sum_{k=1}^K \|\hat{p}_{b,k} - p_{b,k}\|_1 \quad (8)$$

The structure preservation loss function  $\mathcal{L}_{\text{structure}}$  maintains geometric consistency among all detected keypoints. It is computed by measuring the difference between pairwise Euclidean distance matrices of the predicted and ground-truth keypoint sets:

$$\mathcal{L}_{\text{structure}} = \frac{1}{B} \sum_{b=1}^B \|D(\hat{P}_b) - D(P_b)\|_1 \quad (9)$$

where  $\hat{P}_b = \{\hat{p}_{b,1}, \dots, \hat{p}_{b,K}\}$ ,  $P_b = \{p_{b,1}, \dots, p_{b,K}\}$ , and  $D(\cdot)$  denotes the pairwise Euclidean distance matrix computed from a set of keypoints.

By combining these three loss functions, the proposed structure-aware loss function enables precise, robust keypoint detection even under challenging conditions such as motion blur, sparse event data.

### D. Event-based density keypoints tracking

The keypoints of an object detected by the event camera often exhibit pronounced structural saliency. When the object is in a state of high-speed motion, these structurally salient keypoints, due to the presence of multi-directional brightness gradients at their locations, induce diverse intensity changes during motion. This leads to continuous triggering of both positive and negative polarity events around them, forming dense intersections of event streams. As a result, the event density around keypoints increases significantly. Therefore, to achieve robust 6-DoF pose tracking of high-speed objects, this work introduces an EKF-based event density tracking algorithm that leverages event density for matching.

In the proposed keypoints tracking method, we begin by analyzing the polarity distribution within the local region of each keypoint. Based on event polarity statistics, each keypoint is classified as a single-polarity point formed primarily by events of single polarity or a mixed-polarity point resulting from intersecting event lines of different polarities. The classification is determined by the condition described as:

$$\frac{\sum T^+}{\sum(T^+ + T^-)} > \eta \quad \text{or} \quad \frac{\sum T^-}{\sum(T^+ + T^-)} > \eta \quad (10)$$

where  $T^+$ ,  $T^-$  denote the total number of positive events and negative events within the local region, respectively, and  $\eta$  is the threshold parameter for classification. If neither condition is satisfied and the number of surrounding events is sufficiently large, the keypoint is classified as a mixed-polarity point.

For mixed-polarity keypoints, we design a polarity-aware and event-density-guided sliding window matching strategy. Specifically, within a local window centered at the  $i$ -th keypoint, for each candidate position  $(x, y)$ , we extract local responses from the positive and negative polarity event time surfaces and compute the matching score as:

$$S_i(x, y) = \sqrt{T_i^+(x, y) \cdot T_i^-(x', y')} + \beta \cdot D_i(x, y) \quad (11)$$

where  $T_i^+(x, y)$  and  $T_i^-(x', y')$  denote the responses at the  $i$ -th keypoint's local patch on the positive and negative polarity event time surfaces, respectively.  $D_i(x, y)$  is the

event density response in the corresponding local region. The parameter  $\beta$  is a weighting factor that enhances the contribution of structurally informative regions. Finally, the candidate location with the highest score is selected as the updated position of the  $i$ -th keypoint:

$$(x_i^*, y_i^*) = \arg \max_{(x,y)} S_i(x, y) \quad (12)$$

For single-polarity keypoints, when event polarity is highly imbalanced, the mixed-polarity scoring may degrade due to missing information. To address this, we design a sliding search method based solely on the single-polarity time surface, using event density response as the matching criterion.

We define the event response within the local region of a keypoint as:

$$D(x, y) = T^\sigma(x, y), \quad \sigma \in \{+, -\} \quad (13)$$

where  $T^\sigma(x, y)$  denotes the response at position  $(x, y)$  on the time surface of polarity  $\sigma$ . The polarity  $\sigma$  is selected based on the dominant event type in the region surrounding the keypoint.

Within the local search region  $\mathcal{R}_i$  centered at the keypoint, we perform a dual-scale density-guided search to enhance robustness and structural discrimination. For each candidate position  $(x, y) \in \mathcal{R}_i$ , we first compute the integrated event response over a large window:

$$S(x, y) = \sum_{(u,v) \in \mathcal{W}_b(x,y)} D(u, v) \quad (14)$$

To ensure that the selected position exhibits prominent local structural features, we apply a local maximum constraint on the event response:

$$D(x, y) = \max_{(u,v) \in \mathcal{W}_s(x,y)} D(u, v) \quad (15)$$

The  $\mathcal{W}_b(x, y)$  and  $\mathcal{W}_s(x, y)$  denote the large and small window neighborhoods centered at position  $(x, y)$ , respectively. The large window  $\mathcal{W}_b(x, y)$  captures integrated event responses over a wider spatial region to enhance robustness against noise and event sparsity, and the small window  $\mathcal{W}_s(x, y)$  restricts the response to a local neighborhood, ensuring that the selected position corresponds to a local maximum. Among all candidate positions that satisfy the above constraint, we select the one with the highest integrated event response as the updated keypoint location:

$$(x_i^*, y_i^*) = \arg \max_{(x,y) \in \mathcal{R}_i} \{S(x, y) \mid D(x, y)\} \quad (16)$$

This strategy achieves robust keypoint localization by focusing on regions with the highest event density and prominent structural features. Even under conditions of high-speed motion or abrupt illumination changes that result in the loss of events of a single polarity, the method maintains accurate and stable keypoint localization.

To enhance the temporal robustness of keypoint tracking, we introduce the Extended Kalman Filter. The position and velocity of each keypoint are modeled as a state vector, which is propagated using a state transition model and updated

---

### Algorithm 1 Density Keypoint Tracking with EKF

---

```

1: Input: Keypoints  $X$ , Local time surfaces  $S$ 
2: Output: Updated keypoints  $\hat{X}$ 
3: for each keypoint  $x_i$  and time surface  $S_j$  do
4:   if  $x_i = \text{None}$  or  $S_j = \text{None}$  then
5:      $\hat{x}_i \leftarrow \text{None}$ ; continue
6:   end if
7:   if  $\text{EKF}_i$  not initialized then
8:      $\text{EKF}_i.\text{init}(x_i)$ 
9:   end if
10:   $(p_x, p_y) \leftarrow \text{EKF}_i.\text{predict}()$ 
11:   $(T_2^+, T_2^-, x_0, y_0) \leftarrow \text{surface\_info}(S_j)$ 
12:   $D \leftarrow \text{Blur}(T_2^+ + T_2^-)$ 
13:  Extract patch around  $x_i$  from  $T_2^+$  and  $T_2^-$ 
14:  Compute polarity ratio
15:  if is single polarity then
16:    patch  $\leftarrow T_2^+ + T_2^-$ 
17:    peak  $\leftarrow \text{FindPeak}(\text{patch})$ 
18:    best_pt  $\leftarrow \text{peak} + (x_0, y_0)$ 
19:  else
20:    Initialize best_score  $\leftarrow -\infty$ 
21:    for each offset  $(dx, dy)$  in window do
22:       $s \leftarrow \sqrt{T_2^+ \cdot T_2^-} + \beta \cdot D$ 
23:      if  $s > \text{best\_score}$  then
24:        best_pt  $\leftarrow (dx + x_0, dy + y_0)$ 
25:        best_score  $\leftarrow s$ 
26:      end if
27:    end for
28:  end if
29:   $\text{EKF}_i.\text{update}(\text{best\_pt})$ 
30:   $x_{upd}, y_{upd} \leftarrow \text{KF.state}$ 
31:   $x \leftarrow \alpha \cdot x_i + (1 - \alpha) \cdot x_{upd}$ 
32:   $y \leftarrow \alpha \cdot y_i + (1 - \alpha) \cdot y_{upd}$ 
33:  if  $(x, y)$  inside window then
34:     $\hat{x}_i \leftarrow (x, y)$ 
35:    Append to trajectory
36:  else
37:     $\hat{x}_i \leftarrow (x_i, y_i)$ 
38:  end if
39: end for
40: return  $\hat{X}$ 

```

---

with the observation matching the current frame. The overall estimation process can be described as:

$$\mathbf{x}_t = \mathbf{x}_t^{\text{pred}} + \mathbf{K}_t \left( \mathbf{z}_t - \mathbf{H}\mathbf{x}_t^{\text{pred}} \right) \quad (17)$$

where  $\mathbf{x}_t$  denotes the estimated state of the keypoint at time  $t$ , and  $\mathbf{x}_t^{\text{pred}}$  is the predicted state based on the previous frame.  $\mathbf{z}_t$  represents the observed keypoint position from the current matching,  $\mathbf{H}$  is the observation matrix, and  $\mathbf{K}_t$  is the Kalman gain. This mechanism allows the keypoint to be guided by reliable matching results when available, and to rely on motion prediction when observations are missing or noisy, thereby effectively suppressing drift and jitter, and improving tracking stability. The overall density keypoint tracking with EKF method is illustrated in Algorithm 1.

## IV. EXPERIMENTS

In this section, we evaluate the performance of our proposed object 6-DoF pose tracking method on both simulated and real event datasets. The experimental protocol is structured as follows: Section IV-A describes the detailed

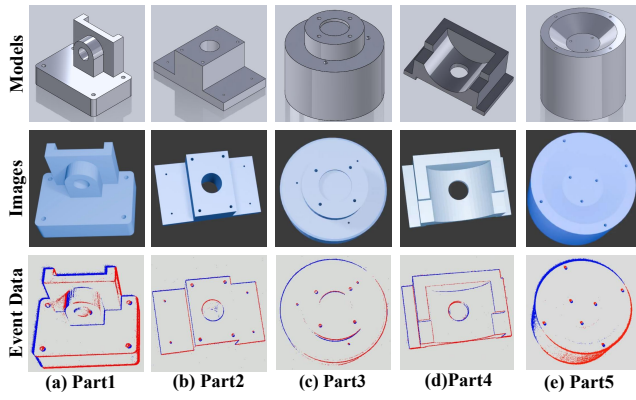


Fig. 3. The models, images, and events data of mechanical parts in simulated event datasets. The top row shows the CAD models of mechanical parts, the middle row shows rendered RGB images of mechanical parts, and the bottom row displays event-accumulated images of those mechanical parts.

implementation of our method; Section IV-B introduces the metrics used for tracking assessment. Then, we present the results of simulated experiments in Section IV-C and real scenario experiments in Section IV-D.

#### A. Implementation Details

In the simulation event experiment, we select several representative mechanical part models, their CAD models, rendered RGB images, and event-accumulated images are shown as illustrated in Fig. 3. The selected models include objects with prominent straight-edge structures as well as those with complex curved contours, to validate the effectiveness of the proposed method in industrial applications.

The steps of the simulation experiment are as follows. Firstly, we employ Blender to render RGB videos of object motion while simultaneously recording the object’s pose variations as ground truth annotations. Subsequently, these videos are converted into event data using the V2E tool [29]. Finally, the error is quantified by comparing the pose tracked by our method with the ground-truth pose computed in Blender.

In the real event experiment, we 3D-print two models derived from the simulation dataset, including a line-based model and a curve-based model, to evaluate the accuracy of our proposed approach. Experiments are performed with an Intel(R) Core(TM) i7-13700 and an NVIDIA GeForce GTX TITAN X GPU (12GB VRAM)

#### B. Evaluation Metrics

To assess the accuracy of the proposed method, we evaluate the 6-DoF object poses using a widely adopted protocol [30] that incorporates the relative pose error.

To evaluate the accuracy of pose tracking, we adopt the relative pose error, denoted as  $R_{rel}$  and  $T_{rel}$ , which measures the discrepancies in rotation and translation between the estimated results and the ground truth. It can be formulated as follows:

$$\mathbf{E}_i = (\mathbf{Q}_i^{-1}\mathbf{Q}_{i+\Delta})^{-1} (\mathbf{P}_i^{-1}\mathbf{P}_{i+\Delta}) \quad (18)$$

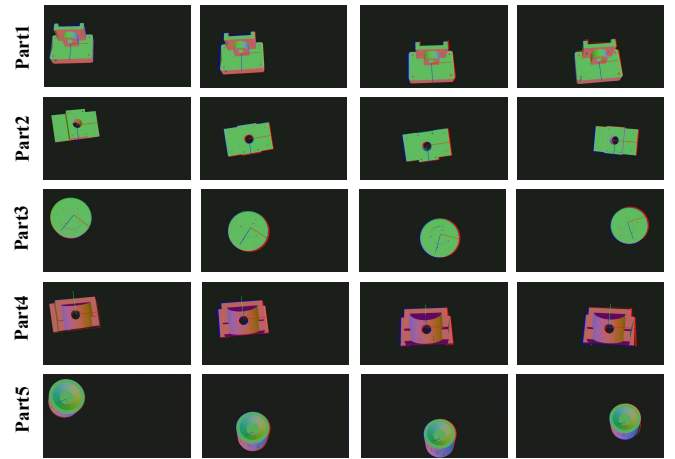


Fig. 4. The simulated event experiment results of high-speed motions, including the rendered model of each object and the 3D coordinate axes representing the predicted object poses.

where  $Q_i$  and  $P_i$  denote the ground truth pose and the estimated pose at the time step  $i$ , respectively.

Subsequently, the relative rotation error  $R_{rel}$  and the relative translation error  $T_{rel}$  are derived from the rotational and translational components of  $E_i$ :

$$R_{rel} = \sqrt{\frac{1}{m} \sum_{i=1}^m \left( \frac{\arccos\left(\frac{\text{trace}(\text{Rot}(E_i)) - 1}{2}\right)}{\Delta t} \right)^2} \quad (19)$$

$$T_{rel} = \sqrt{\frac{1}{m} \sum_{i=1}^m \left( \frac{\|\text{trans}(E_i)\|}{\Delta t} \right)^2} \quad (20)$$

These metrics quantify the relative motion between the ground truth and the estimated trajectory over the same time interval, thereby providing a measure of the local pose drift.

#### C. Simulated Event Experiment

Currently, there is still a lack of event camera datasets specifically designed for mechanical parts 6-DoF pose estimation tasks. Although some existing works have constructed relevant datasets, these have not been publicly released, thus hindering open replication and fair comparison. To address this limitation, we introduce an event-based



Fig. 5. Experimental scenarios of real events.

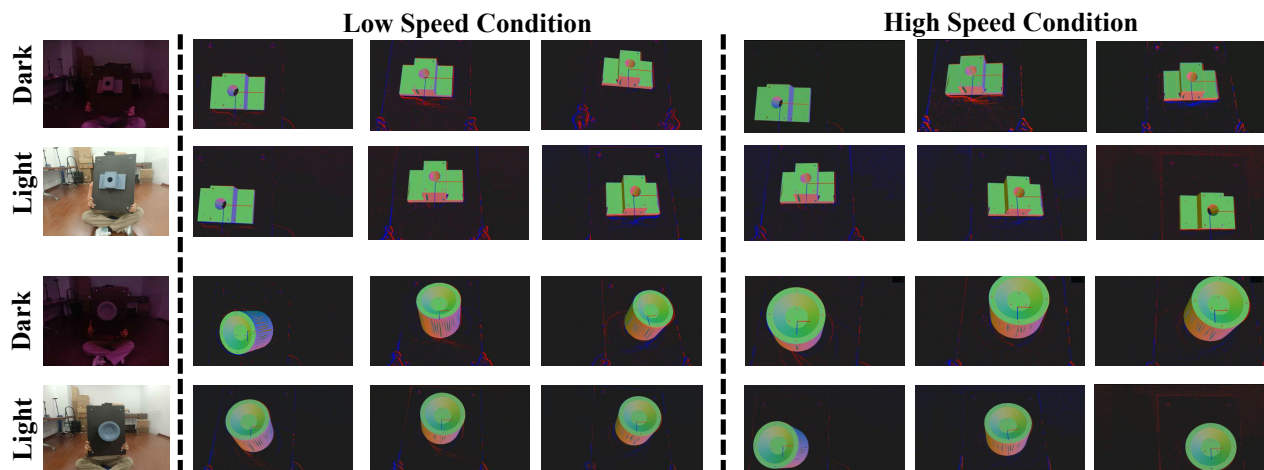


Fig. 6. Pose tracking results of part2 and part5 of real event experiments, including the rendered model of each object and the 3D coordinate axes representing the predicted object poses.

dataset of irregularly moving mechanical parts, comprising both simulated and real events, to facilitate accurate evaluation of the proposed methods in object 6-DoF pose tracking.

In simulated experience, we select five types of mechanical parts with distinct linear and curved features as test objects. Each object is assigned different motion trajectories and velocities, to fully evaluate the adaptability and robustness of the algorithm under various motion patterns. The motion of objects is captured by a stationary monocular camera at a sampling rate of 60 Hz and a resolution of  $1280 \times 720$  pixels, which is consistent with the resolution of the Prophesee EVK4 event camera. Subsequently, the video captured by the monocular camera is converted into event data using the V2E tool.

TABLE I

COMPARISON OF RELATIVE ROTATION ERROR AND RELATIVE TRANSLATION ERROR ON SIMULATED EVENT DATASETS [ $R_{rel}$ : DEG/S,  $T_{rel}$ : CM/S].

Method	Condition		EDOPT[17]		Ours	
	linear velocity	angular velocity	$R_{rel}$	$T_{rel}$	$R_{rel}$	$T_{rel}$
Part1	2.92 m/s	8.27 deg/s	0.81	4.91	<b>0.48</b>	<b>2.97</b>
Part1	5.84 m/s	16.54 deg/s	0.88	2.38	<b>0.28</b>	<b>1.77</b>
Part2	3.51 m/s	31.64 deg/s	1.51	2.00	<b>0.32</b>	<b>1.91</b>
Part2	7.02 m/s	63.28 deg/s	1.04	4.93	<b>0.37</b>	<b>1.86</b>
Part3	2.80 m/s	48.81 deg/s	7.55	<b>3.23</b>	<b>0.69</b>	4.36
Part3	5.59 m/s	97.62 deg/s	7.88	<b>3.32</b>	<b>0.43</b>	4.70
Part4	2.47 m/s	23.82 deg/s	4.47	6.20	<b>0.80</b>	<b>3.26</b>
Part4	4.94 m/s	47.64 deg/s	5.12	4.22	<b>0.67</b>	<b>3.71</b>
Part5	4.12 m/s	29.85 deg/s	7.86	5.40	<b>0.62</b>	<b>1.73</b>
Part5	8.24 m/s	59.70 deg/s	6.71	5.55	<b>0.63</b>	<b>1.62</b>

The experimental results of object 6-DoF pose tracking under simulated event conditions are presented in Fig. 4. This figure presents the 6-DoF pose tracking results of objects in simulated event data from multiple viewpoints, where the green, red, and blue axes indicate the predicted 3D rotation and translation. Meanwhile, the experimental results reported in Table I, compared to the state-of-the-art method, further validate the effectiveness of our proposed approach.

TABLE II

COMPARISON OF RELATIVE ROTATION ERROR AND RELATIVE TRANSLATION ERROR ON THE SELF-COLLECTED REAL EVENT DATASETS [ $R_{rel}$ : DEG/S,  $T_{rel}$ : CM/S].

Method	Condition		EDOPT[17]		Ours	
	Environment	Speed	$R_{rel}$	$T_{rel}$	$R_{rel}$	$T_{rel}$
Part2	Light	Slow	1.21	<b>4.17</b>	<b>0.51</b>	4.21
Part2	Dark	Slow	1.77	4.99	<b>1.03</b>	<b>4.44</b>
Part2	Light	Fast	0.81	7.21	<b>0.44</b>	<b>5.06</b>
Part2	Dark	Fast	1.66	6.98	<b>0.49</b>	<b>5.05</b>
Part5	Light	Slow	7.32	5.66	<b>0.75</b>	<b>3.32</b>
Part5	Dark	Slow	5.11	6.09	<b>0.54</b>	<b>4.76</b>
Part5	Light	Fast	7.02	9.11	<b>1.41</b>	<b>3.40</b>
Part5	Dark	Fast	6.89	8.94	<b>0.46</b>	<b>5.94</b>

The results indicate that our method achieves robust 6-DoF pose tracking, maintaining stable and reliable object tracking performance even in complex and dynamic environments.

#### D. Real Event Experiment

In this section, we conduct practical experimental verification utilizing self-collected real event datasets. The overall experimental scene layout is illustrated in Fig. 5. In the experimental setup, we use the Prophesee EVK4 event camera ( $1280 \times 720$  pixels) to capture event data from 6-DoF object pose changes, while the NOKOV motion capture system is used to acquire the corresponding 6-DoF ground truth pose.

Both devices have been pre-calibrated to ensure spatial and temporal alignment. Subsequently, we select representative target objects with linear contour attributes and curved-surface contour attributes to conduct pose tracking experiments, followed by a systematic comparative analysis of the results. The event camera is placed at a fixed distance from the objects, and we move the objects by hand while the camera captures their motion and records the corresponding events. The 3D models of these objects are known in advance. Markers are attached to these objects to obtain the ground truth of their poses and trajectories using the NOKOV motion capture system.

In order to demonstrate the performance of our method in real-world scenarios, we conduct multiple sets of experiments, and the results are shown in Table II. The object tracking results using the event camera are shown in Fig. 6. The results indicate that the proposed method can achieve robust pose tracking in high-speed motion scenarios.

## V. CONCLUSIONS

This paper proposes a keypoint-based approach for dynamic object 6-DoF pose tracking using event camera. The proposed method achieves stable 6-DoF pose tracking for a wide range of objects, including those with curved geometries, without requiring a predefined initialization. We present a lightweight network for efficient event-based keypoints detection and propose a density-based tracking algorithm with an extended Kalman filter to achieve robust 6-DoF pose tracking under highly dynamic conditions. Extensive experiments demonstrate that the proposed method outperforms state-of-the-art approaches on both synthetic and real-world datasets. However, the proposed method requires known CAD models for object pose tracking. Future work will focus on developing approaches that enable pose tracking for a wider range of objects without relying on CAD models.

## REFERENCES

- [1] G. Du, K. Wang, S. Lian, and K. Zhao, "Vision-based robotic grasping from object localization, object pose estimation to grasp estimation for parallel grippers: a review," *Artificial Intelligence Review*, vol. 54, no. 3, pp. 1677–1734, 2021.
- [2] S. Stevšić, S. Christen, and O. Hilliges, "Learning to assemble: Estimating 6d poses for robotic object-object manipulation," *IEEE Robotics and Automation Letters*, vol. 5, no. 2, pp. 1159–1166, 2020.
- [3] C. Chen, T. Wang, D. Li, and J. Hong, "Repetitive assembly action recognition based on object detection and pose estimation," *Journal of Manufacturing Systems*, vol. 55, pp. 325–333, 2020.
- [4] D. Maji, S. Nagori, M. Mathew, and D. Poddar, "Yolo-6d-pose: Enhancing yolo for single-stage monocular multi-object 6d pose estimation," in *2024 International Conference on 3D Vision (3DV)*. IEEE, 2024, pp. 1616–1625.
- [5] T.-T. Do, T. Pham, M. Cai, and I. Reid, "Real-time monocular object instance 6d pose estimation," in *British Machine Vision Conference 2018*. British Machine Vision Association, 2018.
- [6] S. Wang, J. Liu, Q. Lu, Z. Liu, Y. Zeng, D. Zhang, and B. Chen, "6d pose estimation for vision-guided robot grasping based on monocular camera," in *2023 6th International Conference on Robotics, Control and Automation Engineering (RCAE)*. IEEE, 2023, pp. 13–17.
- [7] T. Pöllabauer, J. Emrich, V. Knauth, and A. Kuijper, "Extending 6d object pose estimators for stereo vision," in *International Conference on Pattern Recognition and Artificial Intelligence*. Springer, 2024, pp. 106–119.
- [8] U. Franke, C. Rabe, H. Badino, and S. Gehrig, "6d-vision: Fusion of stereo and motion for robust environment perception," in *Joint Pattern Recognition Symposium*. Springer, 2005, pp. 216–223.
- [9] J. Ma and J. W. Burdick, "A probabilistic framework for stereo-vision based 3d object search with 6d pose estimation," in *2010 IEEE International Conference on Robotics and Automation*. IEEE, 2010, pp. 2036–2042.
- [10] Y. Hu, S. Speierer, W. Jakob, P. Fua, and M. Salzmann, "Wide-depth-range 6d object pose estimation in space," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 15 870–15 879.
- [11] R. Kaskman, S. Zakharov, I. Shugurov, and S. Ilic, "Homebreweddb: Rgb-d dataset for 6d pose estimation of 3d objects," in *Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops*, 2019, pp. 0–0.
- [12] H. Kim, S. Leutenegger, and A. J. Davison, "Real-time 3d reconstruction and 6-dof tracking with an event camera," in *European conference on computer vision*. Springer, 2016, pp. 349–364.
- [13] G. Gallego, T. Delbrück, G. Orchard, C. Bartolozzi, B. Taba, A. Censi, S. Leutenegger, A. J. Davison, J. Conrad, K. Daniilidis *et al.*, "Event-based vision: A survey," *IEEE transactions on pattern analysis and machine intelligence*, vol. 44, no. 1, pp. 154–180, 2020.
- [14] Z. Wang, Q. Song, Y. Peng, and W. Bai, "Cs3d: An efficient facial expression recognition via event vision," *arXiv preprint arXiv:2512.09592*, 2025.
- [15] Z. Liu, B. Guan, Y. Shang, Q. Yu, and L. Kneip, "Line-based 6-dof object pose estimation and tracking with an event camera," *IEEE Transactions on Image Processing*, 2024.
- [16] Z. Liu, B. Guan, Y. Shang, Y. Bian, P. Sun, and Q. Yu, "Stereo event-based, 6-dof pose tracking for uncooperative spacecraft," *IEEE Transactions on Geoscience and Remote Sensing*, 2025.
- [17] A. Glover, L. Gava, Z. Li, and C. Bartolozzi, "Edopt: Event-camera 6-dof dynamic object pose tracking," in *2024 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2024, pp. 18 200–18 206.
- [18] Z. He, Q. Li, X. Zhao, J. Wang, H. Shen, S. Zhang, and J. Tan, "Contourpose: Monocular 6-d pose estimation method for reflective textureless metal parts," *IEEE Transactions on Robotics*, vol. 39, no. 5, pp. 4037–4050, 2023.
- [19] H. Chen, P. Wang, F. Wang, W. Tian, L. Xiong, and H. Li, "Epropnp: Generalized end-to-end probabilistic perspective-n-points for monocular object pose estimation," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022, pp. 2781–2790.
- [20] D. Reverter Valeiras, G. Orchard, S.-H. Ieng, and R. B. Benosman, "Neuromorphic event-based 3d pose estimation," *Frontiers in neuroscience*, vol. 9, p. 522, 2016.
- [21] D. Reverter Valeiras, S. Kime, S.-H. Ieng, and R. B. Benosman, "An event-based solution to the perspective-n-point problem," *Frontiers in neuroscience*, vol. 10, p. 208, 2016.
- [22] A. Rathinam, H. Qadadri, and D. Aouada, "Spades: A realistic spacecraft pose estimation dataset using event sensing," in *2024 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2024, pp. 11 760–11 766.
- [23] W. Yishi, M. Maestrini, Z. Zexu, M. Massari, and P. Di Lizia, "Cross-modal fusion of monocular images and neuromorphic streams for 6d pose estimation of non-cooperative targets," *Aerospace Science and Technology*, p. 110338, 2025.
- [24] A. G. Howard, M. Zhu, B. Chen, D. Kalenichenko, W. Wang, T. Weyand, M. Andreetto, and H. Adam, "Mobilenets: Efficient convolutional neural networks for mobile vision applications," *arXiv preprint arXiv:1704.04861*, 2017.
- [25] Q. Wang, B. Wu, P. Zhu, P. Li, W. Zuo, and Q. Hu, "Eca-net: Efficient channel attention for deep convolutional neural networks," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 11 534–11 542.
- [26] S. Mehta and M. Rastegari, "Mobilevit: light-weight, general-purpose, and mobile-friendly vision transformer," *arXiv preprint arXiv:2110.02178*, 2021.
- [27] B. Chakravarthi, A. A. Verma, K. Daniilidis, C. Fermuller, and Y. Yang, "Recent event camera innovations: A survey," in *European Conference on Computer Vision*. Springer, 2025, pp. 342–376.
- [28] J. Chen, R. Chen, W. Wang, J. Cheng, L. Zhang, and L. Chen, "Tinyu-net: Lighter yet better u-net with cascaded multi-receptive fields," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2024, pp. 626–635.
- [29] Y. Hu, S.-C. Liu, and T. Delbruck, "v2e: From video frames to realistic dvs events," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2021, pp. 1312–1321.
- [30] J. Sturm, N. Engelhard, F. Endres, W. Burgard, and D. Cremers, "A benchmark for the evaluation of rgb-d slam systems," in *2012 IEEE/RSJ international conference on intelligent robots and systems*. IEEE, 2012, pp. 573–580.