

0-MMS: Zero-Shot Multi-Motion Segmentation With A Monocular Event Camera

Chethan M. Parameshwara, Nitin J. Sanket, Chahat Deep Singh, Cornelia Fermüller, Yiannis Aloimonos

Abstract—Segmentation of moving objects in dynamic scenes is a key process in scene understanding for navigation tasks. Classical cameras suffer from motion blur in such scenarios rendering them effete. On the contrary, event cameras, because of their high temporal resolution and lack of motion blur, are tailor-made for this problem. We present an approach for monocular multi-motion segmentation, which combines bottom-up feature tracking and top-down motion compensation into a unified pipeline, which is the first of its kind to our knowledge. Using the events within a time-interval, our method segments the scene into multiple motions by splitting and merging. We further speed up our method by using the concept of motion propagation and cluster keyframes.

The approach was successfully evaluated on both challenging real-world and synthetic scenarios from the EV-IMO, EED, and MOD datasets and outperformed the state-of-the-art detection rate by 12%, achieving a new state-of-the-art average detection rate of 81.06%, 94.2% and 82.35% on the aforementioned datasets. To enable further research and systematic evaluation of multi-motion segmentation, we present and open-source a new dataset/benchmark called MOD++, which includes challenging sequences and extensive data stratification in-terms of camera and object motion, velocity magnitudes, direction, and rotational speeds.

SUPPLEMENTARY MATERIAL

The accompanying video and dataset are available at prg.cs.umd.edu/0-MMS.html.

I. INTRODUCTION

Navigation is a fundamental competence of life with visual motion estimation as its beating heart [1]–[3]. Even though motion estimation has seen a tremendous advancement in the last few decades, dynamic object motion is usually addressed by outlier rejection schemes as a part of the mature structure from motion and SLAM pipelines [4]. Though, this can provide an initial dynamic object segmentation, further processing for each segment relies on some prior information (commonly appearance/structure/recognition).

To exacerbate the scenario further, classical imaging cameras often fail in dynamic scenarios (moving objects) due to motion blur and low light scenarios. To this end, drawing inspiration from nature, neuromorphic engineers developed a sensor called *Dynamic Vision Sensor* (DVS) [5] which records the asynchronous temporal changes in the scene in the form of a stream of events, rather than the conventional image frames. This gives an unparalleled advantage in-terms of temporal resolution, low latency, and low band-width signals. Such event data is tailor-made for motion segmentation because of the disparity in event density at object boundaries.

This work was supported in parts by the Office of Naval Research and the National Science Foundation under Grants N00014-17-1-2622 and BCS 1824198 respectively and by Samsung Electronics. *Corresponding Author: Chethan M. Parameshwara.*

All authors are associated with the Perception and Robotics Group, University of Maryland, College Park. Emails: {cmparam9, nitin, chahat, fer, yiannis}@umiacs.umd.edu

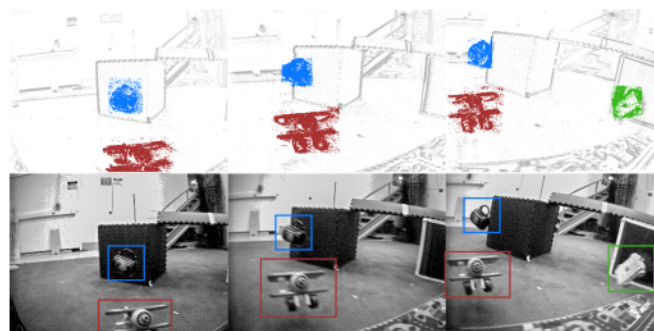


Fig. 1. Multi-Motion Segmentation with a monocular event camera on an EV-IMO dataset sequence. Top Row: The event frames are color-coded by cluster membership. The corresponding grayscale frames are shown in the bottom row. Bounding boxes on the images are color coded with respect to the objects for reference. *Note that grayscale images are not used for computation and are provided for visualization purposes only. All the images in this paper are best viewed in color.*

In this paper, we present a method to detect moving objects by inferring their motion using a monocular event camera; we call this *multi-motion segmentation*. We formally define the problem statement and our contributions next.

A. Problem Formulation and Contributions

We address the following question: *How do you cluster the scene into background and Independently Moving Objects (IMOs) based on motion using data from a moving monocular event camera?*

Given an event volume \mathcal{E} , we find and cluster the events based on 2D motion. We over-segment the scene with the help of feature tracks and then merge clusters based on the motion models and a contrast score. Each cluster is represented by a four parameter motion model (denoting the similarity transformation/warp) $\Theta = \{\Theta_x, \Theta_y, \Theta_z, \Theta_\theta\}$ which represents the 2D translation (Θ_x, Θ_y) , divergence and in-plane rotation, respectively [6]. To speed up computation, we propagate these motion models until a cluster keyframe is invoked. A summary of our contributions is given below (Sample outputs are shown in Fig. 1):

- A novel cluster splitting and merging approach for monocular event-based multi-motion segmentation without prior knowledge of scene geometry (zero-shot) and a number of objects.
- New open-source multi-motion segmentation dataset and benchmark MOD++ including extensive motion stratification and including challenging collision/exploding sequences.
- Speeding up computation using motion propagation and introduction of cluster keyframes.

B. Related Work

There has been extensive progress in the field of event-based motion segmentation in the past decade for

different scenarios at variable scene complexity. Earlier works focused on the case of a static camera, where the events are generated by the moving objects, and a simple clustering scheme can provide motion segmentation [7]–[10]. The next mark up in complexity is the case of a moving camera, where event alignment is computed for the whole scene [6], [11], [12] (also called sharpness or contrast measure [13], [14]) and the parts of the scene which are misaligned give the motion segmentation of IMOs [6] using a simple thresholding algorithm. The results were further improved by [15], who used an Expectation-Maximization scheme to obtain better segmentation. Our work is closely related to [15] with the same underlying philosophy: using motion compensation for clustering but adds robustness in long term segmentation using feature tracking and cluster splitting and merging. We also introduce motion propagation and the concept of cluster keylices to speed-up the entire procedure.

Finally, a few approaches used machine learning [16] to learn object contours and border-ownership information via a structured random forest, which they demonstrated for segmentation. [17], [18] demonstrated a combination of supervised and unsupervised CNN learning using deblurring/event-alignment in the cost function, and [19] designed a graph convolutional neural network for supervised motion segmentation, that uses as input event volumes over extended time periods.

II. PRELIMINARIES

A. Data From An Event Camera

A traditional camera records frames at a fixed frame rate by integrating the number of photons for the chosen shutter time for all pixels *synchronously*. In contrast, an event camera only records the polarity of logarithmic brightness changes *asynchronously* at each pixel. If the brightness at time t of a pixel at location \mathbf{x} is given by $I_{t,\mathbf{x}}$ an event is triggered when $\|\log(I_{t+\delta t,\mathbf{x}}) - \log(I_{t,\mathbf{x}})\|_1 \geq \tau$. Here, δt is a small time increment and τ is a threshold which will determine the trigger of an event (τ is set at the driver level as a combination of multiple parameters). Each event outputs the following data: $\mathbf{e} = \{\mathbf{x}, t, p\}$, where $p = \pm 1$ denotes the sign of the brightness change. We'll denote events in a spatio-temporal window as $\mathcal{E}_t = \{e_i\}_{i=1}^N$ (N is the number of events) and we'll refer to \mathcal{E} as event slice/stream/cloud/volume.

B. Model Fitting For Contrast Maximization

Processing event cloud is generally computationally very expensive and to speed up the processing we use a projection function. The projection of \mathcal{E} leads to a "blurry" image, and a number of methods for measuring this blurriness to achieve event-cloud alignment (also called contrast maximization or motion compensation or deblurring) have been developed [13], [14], [20]. The output of the alignment is an event frame denoted as \mathcal{E}_t . In particular, we utilize the method proposed in [6] to estimate model parameters Θ to maximize the alignment \mathcal{E} by minimizing temporal gradients $\nabla \mathcal{T}$. Here $\bar{\mathcal{T}}(\mathcal{E}) = \mathbb{E}(t_{\mathbf{x}} - t_0)$, \mathbb{E} is the expectation/averaging operator, $t_{\mathbf{x}}$ denotes the time value at location \mathbf{x} , and t_0 is the initial time of the temporal window. Formally, we solve the following optimization problem: $\operatorname{argmin}_{\Theta} \|\nabla \mathcal{T}\|_2$, where ∇ denotes the spatial gradient operator. Note that the projection

Algorithm 1: Splitting and Merging

Data: Tracklets \mathbb{T}_t , Event Stream \mathcal{E}_t , Num. oversegments K
Result: Clusters $\{\mathcal{C}_t\}$, Cluster Models $\{\Theta_t\}$, Segmentation Masks $\{\mathcal{S}\}$
 Splitting;
 $\{\tilde{\mathcal{C}}_t\} = k\text{-Means}(\mathbb{T}_t, K)$;
 Merging;
 $\{\tilde{\Theta}_t\} = \text{ClusterModelFit}(\mathcal{E}_t, \{\tilde{\mathcal{C}}_t\})$;
while *Stopping Criterion and All Clusters Visited do*
 if $\mathcal{C}_{k,j} \mathcal{D}_{k,j}^{-1} > \zeta$ (*Merging Criterion*) **then**
 $\{\tilde{\mathcal{C}}_t\}, \{\tilde{\Theta}_t\} = \text{MergeClusters}(\mathcal{E}, \{\tilde{\mathcal{C}}_t\}, \{\tilde{\Theta}_t\})$
 Updated Clusters and Motion Models;
 end
end
 $\{\mathcal{C}_t\} = \{\tilde{\mathcal{C}}_t\}$ Final Clusters;
 $\{\Theta_t\} = \{\tilde{\Theta}_t\}$ Final Models;
 $\mathcal{S} = \text{ConvexHull}(\{\mathcal{C}_t\})$ Final Dense Segmentation;

function denoted by \mathcal{W} is also called a warping function and refers to the creation of \mathcal{E}_t using parameters Θ .

C. Tracklets Using Point Tracker

We rely on obtaining tracklets (feature tracks across multiple event frames) as an input to multi-motion segmentation. Over the past few years, robust feature extraction and tracking approaches for event data have been proposed [12], [21]–[24], however most of the robust methods are relatively slow or not open-source or use conventional intensity images. Hence, we adapt SuperPoint [25] (previously used on grayscale images) for extracting tracklets \mathbb{T}_t from a set of consecutive N event frames $\{\mathcal{E}_{t+\Delta t \times i} | i \in [0, N-1]\}$. We found the SuperPoint tracker to be robust and generalizable over a wide range of scenarios without any fine tuning. The SuperPoint tracker runs in the backend (we call this Tracker backend) on a First-In First-Out (FIFO) buffer of consecutive N event slices.

III. PROPOSED APPROACH

A. Overview

The proposed solution comprises of two steps: 1. *Split and Merge* summarized in Algorithm 1 and illustrated in Fig. 2. 2. *Motion Propagation and Cluster keylices* summarized in Algorithm 2.

B. Notations

Before we explain the algorithm, we need to define the notations used formally. Tracklets \mathbb{T}_t^i , event stream \mathcal{E}_t^i , clusters \mathcal{C}_t^i , cluster models Θ_t^i at time t for the i^{th} cluster. If a superscript is omitted, the quantity represents the entire data available at time t . Contrast function $\mathcal{C}_{k,j}$ and distance function $\mathcal{D}_{k,j}$ are defined between two clusters k and j . The user chosen thresholds are ζ (merging criterion), λ (stopping criterion), α (scene contrast measure), and χ (propagation threshold).

C. Split And Merge

Splitting: Here, the tracklets from the backend are clustered into k clusters ($k \gg \text{Num. of objects}$) using k -Means clustering for its simplicity and speed (similar to [26], [27]). If a prior on the number of objects or a bound is known, it

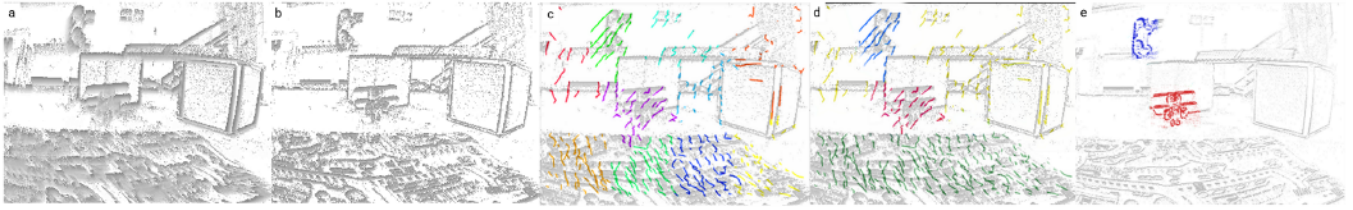


Fig. 2. Overview of the proposed pipeline on a sequence from EV-IMO dataset (a) Projection of the raw event cloud \mathcal{E}_t without motion compensation, (b) Projection of event cloud after global motion compensation (\mathcal{E}_t), (c) Sparse tracklets \mathbb{F}_t extracted on compensated event cloud, (d) Merged feature clusters based on contrast and distance metrics ($\{\mathcal{C}_t\}$), (e) Output of the pipeline is the cluster of events. The cluster membership is color coded where gray color indicating background cluster.

Algorithm 2: Motion Propagation And Cluster Keyslices

Data: Tracklets \mathbb{T}_t , Event Stream \mathcal{E}_t , Clusters $\{\mathcal{C}_{t-1}\}$, Cluster Models $\{\Theta_{t-1}\}$
Result: Clusters $\{\mathcal{C}_t\}$, Cluster Models $\{\Theta_t\}$
foreach Cluster i do
 $\mathcal{E}_t^i = \mathcal{W}(\mathcal{E}_t^i, \mathcal{C}_{t-1}^i, \Theta_{t-1}^i)$ Cluster Event frame;
end
if $\mathbb{E}(\mathcal{C}(\mathcal{E}_t^i) \forall i) > \alpha$ Scene Contrast Measure;
 then
 foreach Cluster i do
 if $\mathcal{C}(\mathcal{E}_t^i) > \chi$ then
 Keep Propagation;
 no new cluster keyslice required;
 else
 Split and Merge on current cluster
 (Algorithm 1);
 New cluster keyslice;
 end
 end
 end
else
 Split and Merge on entire scene (Algorithm 1);
 New cluster keyslice for all clusters (scene);
end

can be trivially incorporated to choose k . We denote these clusters as \mathcal{C}_t .

Merging: Since the splitting method oversegments the scene, we need to merge the clusters to obtain motion segmentation for Independently Moving Objects (IMOs) and the background. The cluster merging is based on a similarity measure that depends on the contrast match (warping a cluster with the model from another cluster and measuring the contrast) and distance between centroids of the clusters. We define contrast and distance functions $\mathcal{C}_{k,j}$ and $\mathcal{D}_{k,j}$ respectively as follows: $\mathcal{C}_{k,j} = \mathbb{E}(\|\text{Var}(\mathcal{E}(\delta\mathcal{E}_j|\Theta_k))\|_1)$ and $\mathcal{D}_{k,j} = \|\mathcal{C}_k - \mathcal{C}_j\|_2$ where k, j are the cluster numbers, $\delta\mathcal{E}_j$ represents the event volume for cluster j and \mathcal{C}_k denotes the centroid of cluster k . This formally entails solving the following optimization problem: $\text{argmax}_j \mathcal{C}_{k,j} \mathcal{D}_{k,j}^{-1}$ which simultaneously maximizes the contrast and minimizes the distance. The larger value of $\mathcal{C}_{k,j} \mathcal{D}_{k,j}^{-1}$ ensures a higher probability that motion model of j^{th} cluster is similar to that of k^{th} cluster (as applying this motion model on $\delta\mathcal{E}_j$ yields maximum contrast). Hence, both clusters belong to one. This step is iteratively performed per cluster (where merging happens with every neighboring cluster using breadth first search) until a stopping criterion has been reached. The

entire process is repeated until all the clusters have been visited. The stopping criterion is explained next.

After each merging operation, we compute the motion model $\Theta_{k,j}$ of the merged clusters by minimizing the temporal gradients $\nabla\mathcal{T}$. Intuitively, when two clusters are merged, the combined motion model captures the average motion of the two clusters thereby slightly increasing the average temporal gradients. Further, whenever a moving object cluster is merged with the background or different IMO cluster the average temporal gradient increases drastically. Hence, a difference in temporal gradient at every step i ($\mathbb{E}(\|\nabla\mathcal{T}_i\|_2)$) with respect to the initial step ($\mathbb{E}(\|\nabla\mathcal{T}_0\|_2)$) is computed. If at any step the difference in temporal gradient is large, we terminate the current merge and continue to the next iteration until all the clusters have been visited. This is mathematically described by $\|\mathbb{E}(\|\nabla\mathcal{T}_i\|_2) - \mathbb{E}(\|\nabla\mathcal{T}_0\|_2)\|_1 \geq \lambda$, where λ is some constant threshold.

After split and merge has been performed, we obtain the final clusters $\{\mathcal{C}_t\}$ ($\{\}$ indicates a set of clusters and each cluster can be indexed with a superscript, i.e., \mathcal{C}_t^i for i^{th} cluster) and motion models per cluster Θ_t^i where i indexes the cluster number. Optionally, we obtain the dense segmentation by taking the convex hull of the sparse feature points in each cluster, which is denoted as \mathcal{S} . Refer to Algorithm 1 for a summary of split and merge methods.

D. Motion Propagation

Optimizing the parameters Θ at every time step to obtain \mathcal{E}_t from \mathcal{E}_t is computationally exorbitant. Inspired by classical tracking pipelines, we propagate motion models from previous to current time slice assuming linear event trajectories. The propagation is achieved using tracklets, and we validate the propagation quality using the following contrast function $\mathcal{C}_t = \mathbb{E}(\|\text{Var}(\mathcal{E}(\delta\mathcal{E}_t|\Theta_{t-1}))\|_1)$. Here, \mathcal{C}_t measures the deviation from the current optimal motion model with respect to the motion model of the previous time slice.

E. Speeding-Up Computation Using Cluster Keyslices

Classical Visual Odometry pipeline utilizes the concept of keyframes to speed-up computation based on certain conditions, this avoids performing bundle adjustment on every frame whilst maintaining good accuracy. We employ a similar strategy of keyframes for event slices, which we call a *keyslice*, for re-clustering based on contrast function \mathcal{C}_t . We apply the contrast function at different levels starting from the entire scene to each clusters separately. Depending on the following two measures the split and merge procedures are performed either per cluster or on the entire scene: 1. Cluster contrast score, 2. Scene contrast score. The cluster contrast

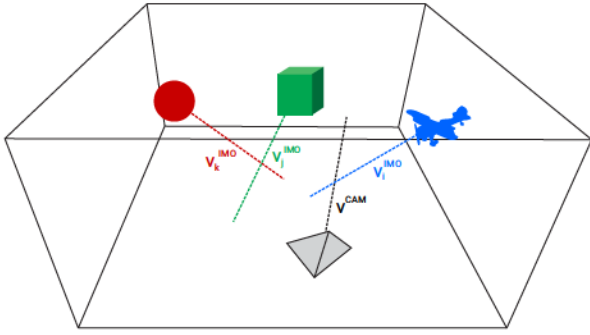


Fig. 3. Velocity vectors v^{CAM} and v_i^{IMO} used for data stratification in the MOD++ dataset.

score is defined by C_t when applied to a single cluster during motion propagation and the average of all cluster scores is called the scene contrast score. Refer to Algorithm 2 for a summary of motion propagation and cluster keyslicing.

IV. MOD++ DATASET/BENCHMARK

Currently, three main datasets exist for IMO segmentation using event cameras, namely, Extreme Event Dataset (EED) [6] and EV-IMO [18] featuring real data and the synthetic Moving Object Dataset (MOD) [17]. However none of the datasets have data stratification based on camera and/or object motion and/or velocities. To this end, we extend the MOD dataset presented in [17] which we call MOD++ to add additional synthetic sequences (Refer to Table I). This data stratification is explained next. Let the instantaneous velocity of the center of mass of the camera and IMOs be denoted as v^{CAM} and v_i^{IMO} , where i is the IMO index (Fig. 3). Now consider the angle and relative-magnitude between two vectors (a, b) denoted by θ and η respectively and defined as follows: $\theta(a, b) = \cos^{-1}\left(\frac{a \cdot b}{\|a\|\|b\|}\right)$ and $\eta(a, b) = \frac{\|a\|}{\|b\|}$. Also, let the instantaneous rotational velocity around its principal axes be denoted by ω where the superscripts and subscripts have the same meaning as that of linear velocities. We classify the sequences as follows: 1. Different linear velocities: Here we set the instantaneous rotational velocity close to zero, i.e., $\|\omega^{\text{CAM}}\| \approx 0$ and $\|\omega_i^{\text{IMO}}\| \approx 0$. We classify the motions based on relative angle and speed. If $\theta(v^{\text{CAM}}, v_i^{\text{IMO}}) \in [0, 45]^\circ \forall i$ we call this sequence small angle. If $\theta(v^{\text{CAM}}, v_i^{\text{IMO}}) \in [60, 120]^\circ \forall i$ we call this sequence medium angle and if $\theta(v^{\text{CAM}}, v_i^{\text{IMO}}) \in [140, 180]^\circ \forall i$ we call this sequence large angle. (Note that we wrap the angles in the range $[0, 180]$ in our case).

If $\eta(v_i^{\text{IMO}}, v_i^{\text{CAM}}) \in [0.5, 3.0] \forall i$ we call this sequence slow. If $\eta(v_i^{\text{IMO}}, v_i^{\text{CAM}}) \in [3.0, 7.0] \forall i$ we call this sequence medium and if $\eta(v_i^{\text{IMO}}, v_i^{\text{CAM}}) \in [7.0, 10.0] \forall i$ we call this sequence fast.

If $\|\omega^{\text{CAM}} - \omega_i^{\text{IMO}}\| \in [0, 5]^\circ s^{-1} \forall i$ we call this sequence slow rotation. If $\|\omega^{\text{CAM}} - \omega_i^{\text{IMO}}\| \in [25, 30]^\circ s^{-1} \forall i$ we call this sequence medium rotation and if $\|\omega^{\text{CAM}} - \omega_i^{\text{IMO}}\| \in [90, 100]^\circ s^{-1} \forall i$ we call this sequence fast rotation.

To make it easy to identify the sequence we use the following naming convention: SeqSEQNUM_ATTR1_ATTRN where SEQNUM is the sequence number which will determine the scene setup (room walls and objects with texture), ATTR1 to ATTRN are modifiers which specify speed and/or rotation classifications. We use the following modifiers: AS, AM, AL for small, medium and large angles, SS, SM, SL for small, medium and large linear speeds, RS, RM, RL

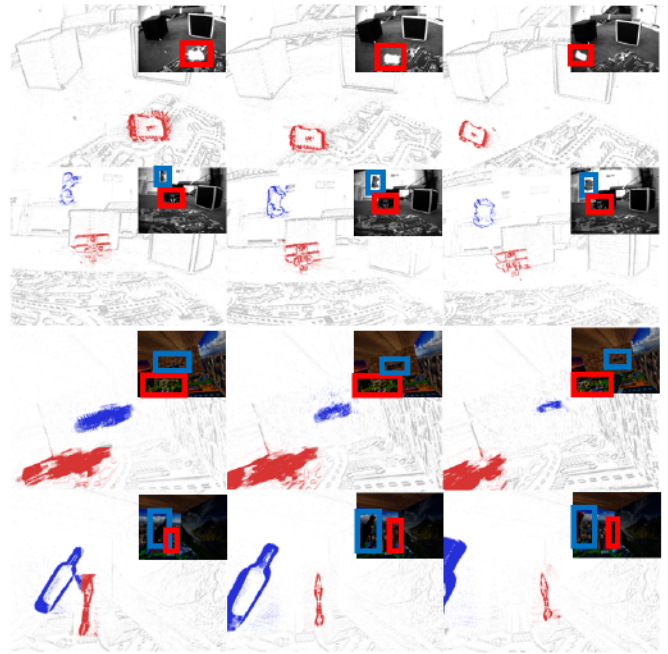


Fig. 4. Qualitative Evaluation of our method on three datasets. Top two rows: EV-IMO dataset, Bottom two rows: MOD dataset. Insets show the corresponding grayscale/RGB images for reference. The cluster membership is color coded where gray color indicates background cluster. Bounding boxes on the images are color coded with respect to the objects for reference.

for small, medium and large rotational speeds. For eg., Seq4_AM_SL_RS would be the fourth sequence with medium angles, large linear speeds and small rotational speeds.

Additionally, we also provide two challenge sequences for researchers to evaluate their algorithm on: Cube and Cup. The Cube sequence is two cubes (a smaller cube on top of a larger cube) falling on the ground and breaking into smaller non-cube pieces. The Cup sequence is a bullet hitting a cup and shattering it into smaller fragments of different shapes.

V. EXPERIMENTS AND RESULTS

We evaluate our approach on publicly available real and synthetic datasets. We demonstrate our approach's performance both qualitatively and quantitatively employing two different metrics based on the availability of groundtruth information.

A. Detection Rate

For datasets which provide timestamped bounding boxes for the objects, we consider the prediction as success when the estimated bounding box fulfills two conditions; (1) it has a overlap of more than at least 50% with the ground-truth bounding box, (2) the area of intersection with the ground-truth box is higher than the intersection with outside area. We can formulate the metric as:

$$\text{Success if } \mathcal{D} \cap \mathcal{G} > 0.5 \text{ and } (\mathcal{D} \cap \mathcal{G}) > (-\mathcal{G} \cap \mathcal{D}) \quad (1)$$

where \mathcal{D} is the predicted mask and \mathcal{G} is the ground truth mask. We evaluate our pipeline's performance on all the datasets using this metric. We obtain the bounding box for our method by obtaining the convex hull on the cluster of events. For comparison purpose we evaluate the performance of [6] using the same metric on all the three datasets. For datasets with more than one sequence, we compute

TABLE I
OVERVIEW OF RELATED DATASETS.

	MOD++	EV-IMO [18]	EED [6]	MOD [17]
Year	2020	2019	2018	2019
Data-type	Simulated	Real	Real	Simulated
Camera	Sim. DAVIS346C	DAVIS346C	DAVIS240B	Sim. DAVIS346C
Data	Events RGB Images @ 1000 Hz 6-DoF Camera Poses @ 1000 Hz 6-DoF IMO Poses @ 1000 Hz IMO Bounding Boxes @ 1000 Hz IMO Masks @ 1000 Hz Optical Flow @ 1000 Hz Depth @ 1000 Hz	Events Grayscale Images @ 40 Hz 6-DoF Camera Poses @ 200 Hz 6-DoF IMO Poses @ 200 Hz IMO Masks @ 40 Hz	Events Grayscale Images @ 20 Hz IMO Bounding Boxes @ 20 Hz	Events RGB Images @ 1000 Hz 6-DoF Camera Poses @ 1000 Hz 6-DoF IMO Poses @ 1000 Hz IMO Bounding Boxes @ 1000 Hz IMO Masks @ 1000 Hz
Poses Ground Truth (Acc.)	Blender [®] Engine (Sub. mm)	12 × Vicon [®] Vantage V8 Cameras (≈ 1mm)	–	Blender [®] Engine (Sub. mm)
IMO Bounding Boxes (Masks) Ground Truth	Blender [®] Engine for both	Scanned 3D Objects projected using Ground Truth Pose	Hand-labelled (–)	Blender [®] Engine for both
Num. Sequences	43	30	5	7
Num. Unique Objects (Max. Number of Objects in frame)	12 (9)	4 (3)	7 (3)	9 (3)
Num. Backgrounds	11	5	5	9
Challenging Scenes	Exploding and Breaking objects	Fast camera motion	Extreme illumination	–
Data Stratification	Velocity Direction Velocity Magnitude IMO Rotation Magnitude	–	–	–

the average of each model’s performance on individual sequences.

B. Intersection Over Union (IoU)

IoU is one the most common and henceforth the most standard measure to evaluate and compare the performance of different segmentation methods. IoU is given by:

$$IoU = (\mathcal{D} \cap \mathcal{G}) / (\mathcal{D} \cup \mathcal{G})$$

Our method outputs a cluster of events which are associated with an object. For the purpose of comparison we convert the sparse mask to a dense mask by assigning all the points lying inside the cluster as the same value.

C. Discussion of Results

Table II reports results on our proposed MOD++ dataset. We pick eight scenarios with different relative (camera and IMO) velocity direction, velocity magnitude and rotation magnitude. We illustrate the merits of split and merge, and motion propagation through extensive ablation studies and compare with previous approaches. Our approach outperforms others by at least ~10%. Executing split and merge at every step offers better accuracy than propagating motion models (evident for more challenging scenarios like slow relative motion, large/small relative velocity direction and/or small rotation magnitude). However, motion propagation offers a speed-up without a significant loss of performance (~1%). Even though simple thresholding [6] and the deep learning-based approach [17] offer better speed-up and Speed×Avg. DR (a metric proposed in [29]), their accuracies are almost 2-3× lower than our approach and is not reliable in challenging scenarios. We leave the speeding-up of our approach using deep learning to enable deployment on mobile robots for future work.

Table III reports the result of our method in comparison with two state-of-the-art IMO detection methods [6], [15] using only a monocular event camera. Our method outperforms the previous methods by up to ~12% detection rate. Specifically, we outperform [6] by a large margin (up to ~32%) on all the three datasets. Our approach is about 2× faster than [15] because of motion propagation while maintaining similar/slightly better accuracy on EED.

Table IV reports the comparison with two deep learning methods for IMO segmentation [18] and [17] using the IoU metric. [17] was trained on the MOD dataset and is tested here on the EV-IMO dataset without any fine-tuning/re-training. We outperform [17] and [18] (which was trained on EV-IMO) on the EV-IMO dataset.

Fig. 4 shows qualitative results of our approach on the two datasets (top two rows show results for the EV-IMO dataset and the last two rows show results for the MOD dataset). Gray areas in the event images show the background cluster and red/blue colored regions show the differently segmented IMOs. The outputs show the robustness of our approach to shape, size and speed of the objects and in-variance with respect to camera motion. Also, note that the objects are sometimes very hard to detect in the corresponding grayscale/RGB frames in Fig. 4 motivating the use of event cameras for IMO detection using motion cues.

We obtain the ground truth IMO by counting ground truth labels with IoU overlap ≤0.2. The graph shows robustness of our approach with increasing number of moving segments. The predicted number of segments closely matches the ground truth. Segmenting solely based on motion with a monocular event-camera is ambiguous in challenging scenarios and results could be improved with incorporation of depth and appearance information in our split and merge step which we believe is the logical next step for future work.

Figs. 5a and 5b show the output of our method for challenging sequences of Gnome shooting and Mug shooting from [28] showing that our method performs well even on real sequences with a large number of objects.

Fig. 6 shows performance of our algorithm with the respect to number of moving segments across time on challenge sequences of MOD++ dataset i.e., Cube and Cup sequence (shown in Figs. 5c and 5d).

Our algorithm runs on a hybrid CPU and GPU system (i7 CPU and NVIDIA Titan Xp GPU). Model fitting and feature extractions are run on GPU in parallel. Even though our core algorithm runs fast, the bottleneck is in the memory transfer to and from the GPU. *The complexity of our approach is linear in the number of clusters, events and frequency of cluster keyframes initiation.* Table II shows the speed of our algorithm in comparison with other approaches in

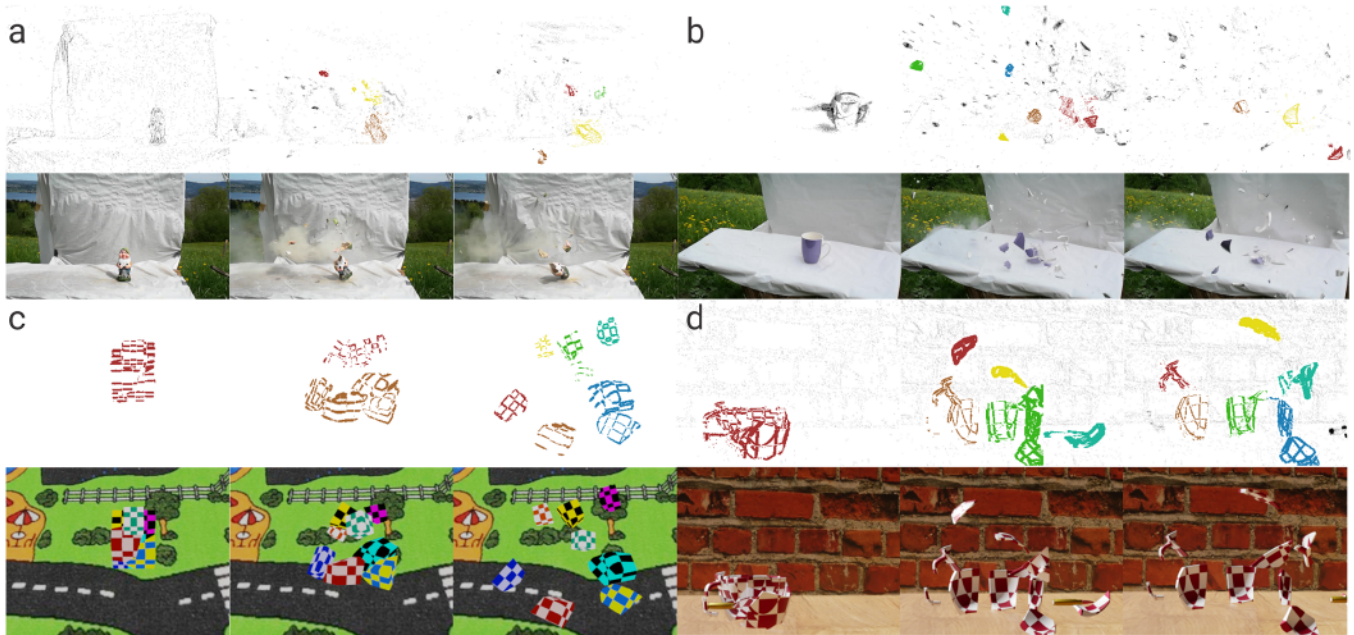


Fig. 5. Multi-Motion Segmentation on the real sequences from [28]. (a) Gnome shooting: Gnome statue getting shot by a bullet, (b) Mug shooting: Mug getting shot by a bullet. Challenge sequences from the proposed MOD++ dataset. (c) Cube breaking into smaller pieces by falling, and (d) Cup getting shot by a bullet. The event frames are colored by cluster membership with gray showing background cluster. Note that the corresponding RGB frames are not used for computation and are provided for visualization purposes only.

TABLE II
COMPARISON WITH STATE-OF-THE-ART USING THE DETECTION RATE FOR DIFFERENT SEQUENCES OF MOD++.

Method	Detection Rate (DR in %) \uparrow									Speed \uparrow (MEv/s)	Speed \times Avg. DR \uparrow
	AS_SM_RS	AM_SM_RS	AL_SM_RS	AM_SS_RS	AM_SL_RS	AL_SS_RS	AL_SS_RM	AL_SS_RL	AL_SS_RS		
Mitrokhin <i>et al.</i> [6]	35.24	32.29	38.12	28.78	43.28	24.56	32.29	39.65	5.41	185.47	
EVDodgeNet [17]	42.25	46.94	53.23	37.81	61.72	46.13	43.50	52.38	10.01	480.42	
<i>k</i> -Means (<i>k</i> =5)	44.89	47.73	49.35	40.17	59.52	42.19	45.71	55.83	1.07	51.54	
<i>k</i> -Means (<i>k</i> =10)	60.36	64.87	59.27	47.73	65.78	48.92	54.74	58.47	1.02	58.66	
<i>k</i> -Means (<i>k</i> =20)	56.1	62.28	58.01	45.25	61.25	44.71	49.37	54.91	0.98	52.90	
Ours (No Propagation)	70.13	73.29	72.58	65.80	85.76	66.24	72.90	79.02	0.83	60.77	
Ours	69.52	73.94	71.27	63.58	84.21	64.93	71.18	78.37	1.16	83.67	

TABLE III
COMPARISON WITH STATE-OF-THE-ART USING DETECTION RATE FOR EED, MOD, EV-IMO DATASETS.

Method	Detection rate for dataset (%) \uparrow			Speed \uparrow (MEv/s)
	EED	MOD	EV-IMO	
Mitrokhin <i>et al.</i> [6]	88.93	70.12	48.79	5.41
Stoffregen <i>et al.</i> [15]	93.17	-	-	0.64* ($N_I=10$)
Ours	94.2	82.35	81.06	1.16

*Results taken directly from [15]

TABLE IV
COMPARISON WITH STATE-OF-THE-ART USING IoU FOR EV-IMO.

Method	IoU
EV-IMO [18]	77.00*
EVDodgeNet [17]	65.76
Ours	80.37

*Results taken directly from [18] in which boxes and wall are used for training.

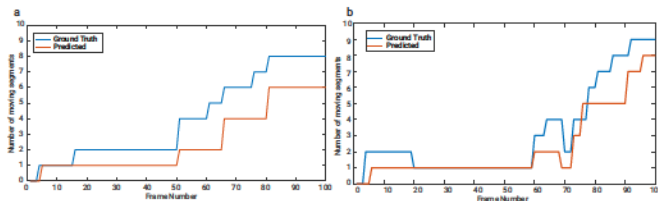


Fig. 6. Number of moving segments vs. frame number (time) on challenge sequences of MOD++: (a) Cube, (b) Cup.

Million Events per second (MEv/s). Our motion propagation and keyslicing provides a speed-up of upto 40% without

compromising accuracy.

VI. CONCLUSIONS AND FUTURE WORK

We presented a method for multi-motion segmentation using data from a monocular event camera. Our approach works by splitting the scene into smaller motions and then iteratively merging them based on a contrast measure. To our knowledge, this is the first approach for monocular independent motion segmentation which combines a bottom-up feature tracking and top-down motion compensation into a unified pipeline. We further speed up our method by using the concept of motion propagation and cluster keyslices.

A comprehensive qualitative and quantitative evaluation is provided on three challenging event motion segmentation datasets, namely, EV-IMO, EED and MOD showcasing the robustness of our approach. Our method outperforms the previous state-of-the-art approaches by upto $\sim 12\%$ detection, thereby achieving the new state-of-the-art on the three aforementioned datasets. To accelerate further research in this area, we present and open-source a new benchmark dataset MOD++ which includes challenging scenes such as cube breaking and a mug getting shot by a bullet along with extensive data stratification in-terms of camera and object motion, velocity magnitudes, direction and rotational speeds. We achieve 73.21% detection rate on MOD++ which is 2 to 3 \times higher than the state-of-the-art methods.

REFERENCES

- [1] Cornelia Fermüller. Navigational preliminaries. *Active Perception*, 1:103–150, 1993.
- [2] Nitin J Sanket, Chahat Deep Singh, Kanishka Ganguly, Cornelia Fermüller, and Yiannis Aloimonos. Gapflyt: Active vision based minimalist structure-less gap detection for quadrotor flight. volume 3, pages 2799–2806. IEEE, 2018.
- [3] Nitin J Sanket, Chahat Deep Singh, Varun Asthana, Cornelia Fermüller, and Yiannis Aloimonos. Morpheys: Variable baseline stereo for quadrotor navigation. *arXiv preprint arXiv:2011.03077*, 2020.
- [4] Huai-Jen Liang, Nitin J Sanket, Cornelia Fermüller, and Yiannis Aloimonos. Salientdso: Bringing attention to direct sparse odometry. *IEEE Transactions on Automation Science and Engineering*, 16(4):1619–1626, 2019.
- [5] Patrick Lichtsteiner, Christoph Posch, and Tobi Delbruck. A 128×128 120 db $15\mu\text{s}$ latency asynchronous temporal contrast vision sensor. *IEEE journal of solid-state circuits*, 43(2):566–576, 2008.
- [6] Anton Mitrokhin, Cornelia Fermüller, Chethan Parameshwara, and Yiannis Aloimonos. Event-based moving object detection and tracking. *IEEE/RSJ Int. Conf. Intelligent Robots and Systems (IROS)*, 2018.
- [7] M. Litzberger, C. Posch, D. Bauer, A. N. Belbachir, P. Schon, B. Kohn, and H. Garn. Embedded vision system for real-time object tracking using an asynchronous transient vision sensor. In *IEEE 12th Digital Signal Processing Workshop and 4th IEEE Signal Processing Education Workshop*, pages 173–178, 2006.
- [8] Alejandro Linares-Barranco, Francisco Gómez-Rodríguez, Vicente Villanueva, Luca Longinotti, and Tobi Delbrück. A usb3. 0 fpga event-based filtering and tracking framework for dynamic vision sensors. In *2015 IEEE International Symposium on Circuits and Systems (ISCAS)*, pages 2417–2420. IEEE, 2015.
- [9] Abhishek Mishra, Rohan Ghosh, Jose C Principe, Nitish V Thakor, and Sunil L Kukreja. A saccade based framework for real-time motion segmentation using event based vision sensors. *Frontiers in neuroscience*, 11:83, 2017.
- [10] Francisco Barranco, Cornelia Fermüller, and Eduardo Ros. Real-time clustering and multi-target tracking using event-based sensors. In *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 5764–5769, 2018.
- [11] Guillermo Gallego and Davide Scaramuzza. Accurate angular velocity estimation with an event camera. *IEEE Robotics and Automation Letters*, 2(2):632–639, 2017.
- [12] A. Z. Zhu, N. Atanasov, and K. Daniilidis. Event-based visual inertial odometry. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5816–5824, July 2017.
- [13] Guillermo Gallego, Mathias Gehrig, and Davide Scaramuzza. Focus is all you need: Loss functions for event-based vision. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 12280–12289, 2019.
- [14] Timo Stoffregen and Lindsay Kleeman. Event cameras, contrast maximization and reward functions: An analysis. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1658–1667, 2020.
- [15] Timo Stoffregen, Guillermo Gallego, Tom Drummond, Lindsay Kleeman, and Davide Scaramuzza. Event-based motion segmentation by motion compensation. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 7244–7253, 2019.
- [16] Francisco Barranco, Cornelia Fermüller, and Yiannis Aloimonos. Bio-inspired motion estimation with event-driven sensors. In *International Work-Conference on Artificial Neural Networks*, pages 309–321. Springer, 2015.
- [17] Nitin J. Sanket, Chethan M. Parameshwara, Chahat Deep Singh, Ashwin V. Kuruttukulam, Cornelia Fermüller, Davide Scaramuzza, and Yiannis Aloimonos. Evdodgenet: Deep dynamic obstacle dodging with event cameras, 2019.
- [18] Anton Mitrokhin, Chengxi Ye, Cornelia Fermüller, Yiannis Aloimonos, and Tobi Delbruck. EV-IMO: motion segmentation dataset and learning pipeline for event cameras. In *IEEE/RSJ Int. Conf. Intelligent Robots and Systems (IROS)*, 2019.
- [19] Anton Mitrokhin, Zhiyuan Hua, Cornelia Fermüller, and Yiannis Aloimonos. Learning visual motion segmentation using event surfaces. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14414–14423, 2020.
- [20] Guillermo Gallego, Henri Rebecq, and Davide Scaramuzza. A unifying contrast maximization framework for event cameras, with applications to motion, depth, and optical flow estimation. *IEEE Conf. Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [22] Daniel Gehrig, Henri Rebecq, Guillermo Gallego, and Davide Scaramuzza. Èklt: Asynchronous photometric feature tracking using events and frames. *International Journal of Computer Vision*, 128(3):601–618, 2020.
- [23] Laurent Dardet, Sio-Hoi Ieng, and Ryad Benosman. Event-based features selection and tracking from intertwined estimation of velocity and generative contours. *arXiv preprint arXiv:1811.07839*, 2018.
- [24] Ignacio Alzugaray and Margarita Chli. Asynchronous multi-hypothesis tracking of features with event cameras. In *2019 International Conference on 3D Vision (3DV)*, pages 269–278. IEEE, 2019.
- [25] Daniel DeTone, Tomasz Malisiewicz, and Andrew Rabinovich. Superpoint: Self-supervised interest point detection and description. In *CVPR Deep Learning for Visual SLAM Workshop*, 2018.
- [26] Qi Wang, Mulin Chen, Feiping Nie, and Xuelong Li. Detecting coherent groups in crowd scenes by multiview clustering. *IEEE transactions on pattern analysis and machine intelligence*, 42(1):46–58, 2018.
- [27] Thomas Brox and Jitendra Malik. Object segmentation by long term analysis of point trajectories. In *European conference on computer vision*, pages 282–295. Springer, 2010.
- [28] Henri Rebecq, René Ranftl, Vladlen Koltun, and Davide Scaramuzza. High speed and high dynamic range video with an event camera. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2019.
- [29] Nitin J Sanket, Chahat Deep Singh, Cornelia Fermüller, and Yiannis Aloimonos. PRGFlow: Benchmarking SWAP-Aware Unified Deep Visual Inertial Odometry. *arXiv preprint arXiv:2006.06753*, 2020.