



**University of  
Zurich**<sup>UZH</sup>

**Zurich Open Repository and  
Archive**

University of Zurich  
University Library  
Strickhofstrasse 39  
CH-8057 Zurich  
[www.zora.uzh.ch](http://www.zora.uzh.ch)

---

Year: 2024

---

## **Dense Continuous-Time Optical Flow from Event Cameras**

Gehrig, Mathias ; Muglikar, Manasi ; Scaramuzza, Davide

DOI: <https://doi.org/10.1109/TPAMI.2024.3361671>

Posted at the Zurich Open Repository and Archive, University of Zurich

ZORA URL: <https://doi.org/10.5167/uzh-257404>

Journal Article

Accepted Version



The following work is licensed under a Creative Commons: Attribution-NonCommercial-NoDerivatives 4.0 International (CC BY-NC-ND 4.0) License.

Originally published at:

Gehrig, Mathias; Muglikar, Manasi; Scaramuzza, Davide (2024). Dense Continuous-Time Optical Flow from Event Cameras. IEEE Transactions on Pattern Analysis and Machine Intelligence:1-12.

DOI: <https://doi.org/10.1109/TPAMI.2024.3361671>

# Dense Continuous-Time Optical Flow from Event Cameras

Mathias Gehrig, Manasi Muglikar, and Davide Scaramuzza

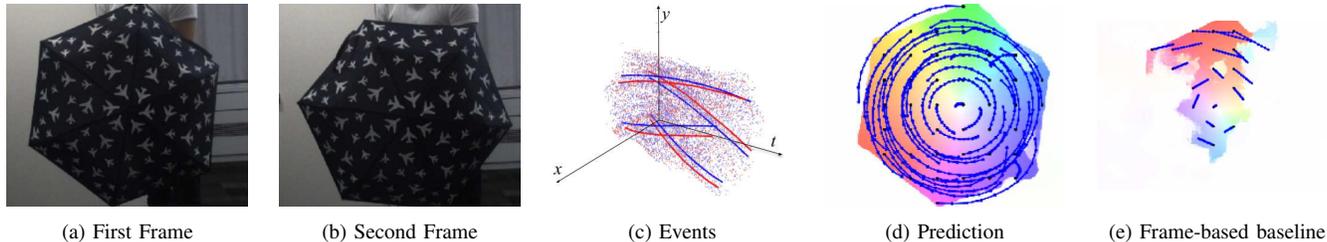


Fig. 1. (d) Our method, trained only on simulated events, predicts dense pixel trajectories in continuous-time on real sequences. The frames in (a) and (b) are shown only for illustrative purposes. In (c), we also visualize selected pixel trajectories in the space-time volume of events. (e) The frame-based baseline [1] attempts to predict linear motion from the frames but falls short due to perceptual ambiguity and non-linear motion in pixel space.

**Abstract**—We present a method for estimating dense continuous-time optical flow from event data. Traditional dense optical flow methods compute the pixel displacement between two images. Due to missing information, these approaches cannot recover the pixel trajectories in the blind time between two images. In this work, we show that it is possible to compute per-pixel, continuous-time optical flow using events from an event camera. Events provide temporally fine-grained information about movement in pixel space due to their asynchronous nature and microsecond response time. We leverage these benefits to predict pixel trajectories densely in continuous time via parameterized Bézier curves. To achieve this, we build a neural network with strong inductive biases for this task: First, we build multiple sequential correlation volumes in time using event data. Second, we use Bézier curves to index these correlation volumes at multiple timestamps along the trajectory. Third, we use the retrieved correlation to update the Bézier curve representations iteratively. Our method can optionally include image pairs to boost performance further. To the best of our knowledge, our model is the first method that can regress dense pixel trajectories from event data. To train and evaluate our model, we introduce a synthetic dataset (MultiFlow) that features moving objects and ground truth trajectories for every pixel. Our quantitative experiments not only suggest that our method successfully predicts pixel trajectories in continuous time but also that it is competitive in the traditional two-view pixel displacement metric on MultiFlow and DSEC-Flow. Open source code and datasets are released to the public.

## MULTIMEDIA MATERIAL

Code and dataset available at:

<https://github.com/uzh-rpg/bflow>

## I. INTRODUCTION

Optical flow estimation is a fundamental low-level vision task that informs about motion in pixel space. It has numerous practical applications in computational photography and videography, video compression, inverse graphics, object tracking, and robotics [2].

The authors are with the Robotics and Perception Group, affiliated with both the Dept. of Informatics of the University of Zurich and the Dept. of Neuroinformatics of the University of Zurich and ETH Zurich, Switzerland.

Traditionally, this problem has been addressed by finding dense correspondences between two frames. Frame-based sensors provide information at a fixed frequency, independent of the dynamics in the scene, and must strike a trade-off between bandwidth and latency: At high speeds, they require a high frame-rate to reduce perceptual latency, but this introduces a significant bandwidth-overhead for downstream systems. Instead, reducing the frame-rate reduces the bandwidth requirements but at the cost of missing important scene dynamics. These shortcomings have inspired recent work to address optical flow estimation with event cameras [3], [4]. In contrast to frame-based sensors, which capture frames at regular intervals, event cameras register per-pixel brightness changes asynchronously and at very high temporal resolution ( $< 1$  millisecond). Additionally, they are robust to motion blur and have a very high dynamic range [5]. Due to these properties, event cameras are ideally suited for optical flow estimation. Nonetheless, these advantages are tied to the fundamentally different data format resulting from event cameras.

Effectively extracting motion information from event cameras is a non-trivial task for which deep neural networks have shown promising results. Until now, neural network approaches regress optical flow at a specific timestamp [3], [4]. Event cameras provide visual cues in continuous-time, which means that predicting discrete pixel displacements ignores much of the device’s potential [4].

In this work, instead, we go a step further and propose a differentiable model that can exploit the rich spatio-temporal nature of event data by regressing continuous-time trajectories for every pixel of the camera. Figure 1 illustrates an example prediction of our method using only events in a short time window. This approach not only enables new applications that can make use of continuous-time pixel trajectories [6], [7] but also generally improves the accuracy of the model, as we will show in our experiments in Sec. V.

Extracting per-pixel continuous-time trajectories from

event data is a challenging problem that requires careful modelling of the problem. To achieve this, we propose the following methodological innovations:

- Instead of predicting pixel displacements, we generalize this concept and estimate control points of Bézier curves for each pixel. As a result, our method can regress the trajectory of each pixel at arbitrary times. It also enables the second innovation:
- We use multiple correlation volumes in time to search for pixel correspondences. We then use the estimated Bézier curves to retrieve correlation features of all correlation volumes simultaneously. This enables the incorporation of motion prior and facilitates the task of finding accurate pixel trajectories.
- The use of image data is optional. We show that our approach works well purely with event data and enable dense, pixel-wise trajectories. Furthermore, we show that a combination of image and event data outperforms single-modality approaches on both real and synthetic datasets.

To train and evaluate our model we also contribute a new synthetic dataset MultiFlow which is inspired by the FlyingChairs dataset [8]. Differently from the FlyingChairs dataset, MultiFlow features event data from various moving objects undergoing continuous similarity transformations as well as dense pixel-trajectory ground truth.

The proposed approach generalizes both the RAFT [9] and E-RAFT [4] architecture by not only giving the option to compute multiple correlation volumes in time but also to predict Bézier curves that include the linear motion model as a special case. The introduction of Bézier curves gives the model the capabilities to estimate non-linear pixel-trajectories but also reduces the end-point-error that is used to assess the accuracy of predicted pixel displacements.

Our experiments suggest that simply giving the model the flexibility to predict Bézier curves reduces the end-point-error on MultiFlow by approximately 67%. Additionally, the introduction of multiple correlation lookups in time further reduces the error metric computed on the whole trajectory by approximately 40%. Finally, we provide quantitative real-world experiments for the traditional optical flow task on DSEC-Flow [4], reducing the end-point-error by 14.5% compared to prior art.

## II. RELATED WORK

### A. Optical Flow

For conciseness, we focus mostly on neural-network-based methods.

1) *Image-based*: The vast majority of neural network-based optical flow method considers the task of estimating dense pixel displacements from a pair of frames [9], [10], [11], [12], [13], [14], [15], [16], [17], [18], [19], [20], [21]. A common component of many highly successful methods are explicit correlation volumes that guide the matching process. This inductive bias enables high performance and data efficiency [9] as well as strong cross-dataset generalization [17].

Multi-Frame optical flow estimation has been mostly explored in the self-supervised learning setting [22], [23], [24] or optimization-based literature [25], [26], [27]. Some neural-network-based approaches use an additional pair of frames to initialize the optical flow prediction [9] or use warped flow as input [28] to a second stage. Recent work proposes to track single pixels from video [29]. By estimating the trajectory only for a single pixel, this method does not utilize spatial information but still achieves remarkable performance. From these multi-frame approaches, Slow Flow [25] is conceptually most related. Slow Flow estimates optical flow  $\mathbf{F}_{1 \rightarrow N}$  given a sequence of  $N$  images at high frame rate (1000 fps) with dense tracking using an optimization-based approach. However, it requires a high-speed camera and complex occlusion handling that incurs a trade-off between drift and accurate prediction at motion boundaries.

2) *Event-based*: Event-based optical flow algorithms can be categorized into four categories: (i) Asynchronous methods [30], [31] using the Lucas-Kanade algorithm [32]. (ii) Plane fitting-based methods that exploit the local plane-like shape of spatio-temporal event streams [33], [34]. (iii) Variational optimization-based approaches [35], [36] that incorporate image data [36] or simultaneously estimate image intensity [35]. (iv) Learning-based approaches, most of which are trained via self-supervision [3], [37], [38], [39], [40]. Supervision is provided either by images [3], [37], [38] or events [39], [40]. Our approach is related to E-RAFT [4], which adapts the RAFT [9] framework to event data to leverage correlation features from cost volumes to estimate dense pixel correspondences for large displacements. In contrast to E-RAFT or RAFT, we predict pixel trajectories using multiple correlation lookups in time while we also show the advantages of combining events and frames.

### B. Continuous Tracking

Continuous-time trajectory estimation of camera poses has been proposed in the context of rolling shutter compensation [41] as well as visual-inertial odometry for event cameras [42]. Instead, we are interested in regressing pixel-trajectories, which is more closely related to high-speed feature tracking for event cameras [31], [43], [44], [45]. Our approach is related to the work of Seok et al. [43] who use quadratic Bézier curves to sparsely track features by maximizing the variance of the image of warped events on local patches. In contrast to our work, this method can only sparsely track features where events are present, cannot incorporate learned priors from data, and does not offer the possibility to include images.

### C. Datasets for Optical Flow

Existing datasets can be categorized as image-based optical flow datasets and event-based optical flow datasets.

1) *Image-based*: The seminal work of FlowNet [21] proposed a large synthetic dataset called "FlyingChairs" to train their CNN. FlyingThings3D [46] introduced a synthetic stereo-video dataset with scene flow ground truth.

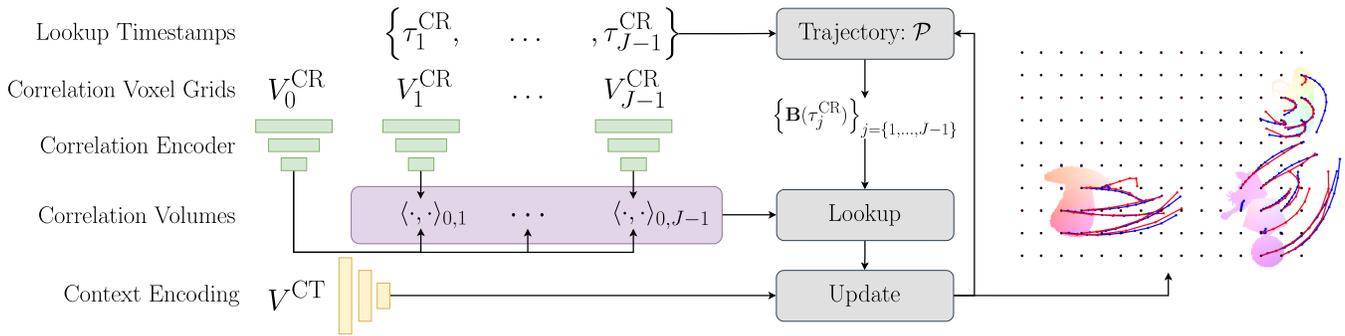


Fig. 2. Overview of the framework. For concise presentation, we illustrate the purely event-based approach. We create multiple correlation voxel grids in a sequence to compute correlation volumes. Bézier curves  $\mathbf{B}$ , parameterized by a set of control points  $\mathcal{P}$ , represent the continuous-time pixel trajectories. Along these trajectories, we index  $J$  correlation volumes at their associated timestamps  $\tau^{\text{CR}}$ . An learned update operator uses the output of this lookup operation to update the Bézier curves. Finally, this loop is repeated in an iterative fashion.

AutoFlow[47] proposes to render large-scale training data with rich data augmentation in 2D for optical flow estimation. These datasets, however, do not provide event data and ground truth pixel trajectories beyond two views.

2) *Event-based*: MVSEC [3] contains 5 outdoor driving sequences and 4 indoor sequences. However, optical flow ground truth from MVSEC suffers from inaccuracies in calibration [4] and only features very small displacements. DSEC-Flow, a more recent dataset, addresses these shortcomings by providing accurate but sparse optical flow ground truth. The main downside of MVSEC and DSEC-Flow is that they both do not include optical flow ground truth for dynamic objects.

In contrast to the aforementioned datasets, our synthetic dataset *MultiFlow* features accurate and dense ground truth for pixel trajectories as well as image and event data. The concept of pixel trajectories is a generalization of pixel displacement and can be used to supervise and evaluate pixel tracking methods.

### III. METHODOLOGY

This paragraph outlines our methodology, setting the stage for the following detailed explanations. To assist in visualizing the process, Figure 2 serves as a helpful guide. Our approach begins by transforming events into a spatio-temporal event representation, crucial for subsequent feature extraction, as elaborated in Section III-B. Following this, our model extracts two key types of features: context and correlation features from the event representations. The intricacies of this process are elaborated in Section III-C. At the heart of our approach lies the computation of correlation volumes. These volumes are critical in assessing the quality of feature matches over time, thereby aiding in the tracking of pixel movements. This aspect is thoroughly explained in Section III-D. Further, we describe our iterative approach to refine pixel motion, which is parameterized by Bézier curves initialized as straight lines. These curves are crucial for indexing the correlation volumes and thus, for guiding the trajectory estimation process. The nuances of this step are covered in Section III-E.

Our implementation reuses many components of the original RAFT architecture [9]. While Figure 2 offers a conceptual summary of our approach, the specific neural network architecture that integrates these elements is depicted in Figure 2 of the supplementary material.

#### A. Problem Definition

Our method is tasked with estimating a function

$$\mathbf{B} : \mathcal{T} \times \mathbb{N}_0 \times \mathbb{N}_0 \rightarrow \mathbb{R}^2 \quad (1)$$

$$(\tau, x, y) \mapsto \mathbf{B}(\tau, x, y) \quad (2)$$

that describes the per-pixel trajectories in time on the image plane.  $\mathcal{T} = \mathbb{R} \cap [0, 1]$  describes the domain of the normalized time  $\tau(t) = (t-t_r)/(t_t-t_r)$ , where  $\tau = 0$  corresponds to the reference time  $t_r$  from which the pixel trajectory starts and  $\tau = 1$  corresponds to the target time  $t_t$  when the pixel trajectory ends. This formulation can be used to find the pixel displacement at any time between  $t_r$  and  $t_t$ .

As will be explained in more detail later, the proposed approach also uses additional events from time  $t_0$  until the reference time  $t_r$  to extract feature maps at different timestamps. However, the pixel trajectory is estimated only from the reference time  $t_r$  to the target time  $t_t$ .

#### B. Input Data Preparation

The proposed approach uses features extracted from event data and optionally a pair of images to further boost performance.

1) *Multi-View Event Representations*: Event cameras have a fundamentally different working principle than frame-based cameras. Instead of acquiring frames, event cameras receive asynchronous events with high temporal resolution. An event  $e_k(t) = (x_k, y_k, t, p_k)$  is a tuple containing information about a pixel  $(x_k, y_k)$  for which a positive or negative brightness change  $p_k$  was registered at time  $t$ . Note that the time  $t$  typically has microsecond resolution, which provides precise temporal information about motion in the scene.

The first step of the feature extraction pipeline is the construction of a discrete spatio-temporal representation from a sequence of events. For our experiments, we choose the voxel grid representation by Zhu et al. [37] due to its simplicity

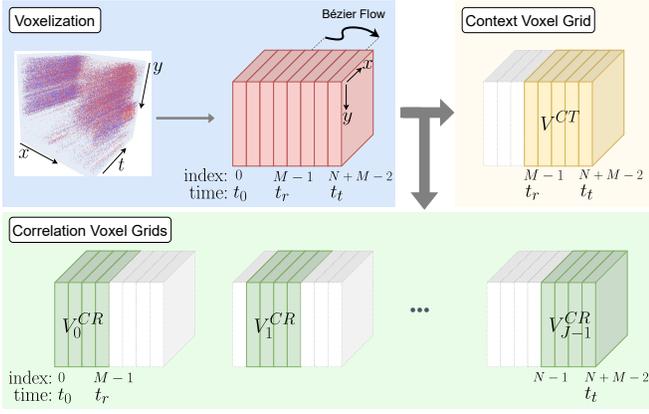


Fig. 3. Voxel Grid Generation. Events are interpolated into a voxel grid representation before being divided into sub-voxel grids that are either used as input to the context encoder or the correlation encoder. The Bézier curves, indicated in the figure, are estimated from the reference time  $t_r$  to the target time  $t_t$  for each pixel. The initial voxel grid channels before the reference time  $t_r$  are used for extracting correlation voxel grids as visualized in the lowest, green row of the figure.

and possibility to extract features in a sliding window. Figure 3 provides an overview of this process.

Given the task of estimating continuous pixel trajectories from  $t_r$ , the reference time, to the target time  $t_t$ , Our method first computes a *base* voxel grid, visualized in red in Figure 3, consisting of  $M + N - 1$  discrete bins along the time dimension via interpolation of event data [37]. The first  $M$  bins are computed from events between initial time  $t_0$  and the reference time  $t_r$ . These events are used to extract the correlation features later in the pipeline. The last  $N$  bins are computed from events during the time window  $[t_t, t_r]$  when the pixel trajectories are estimated. These  $N$  bins are required to extract the context features. The bin at the reference timestamp  $t_r$  will be reused by both correlation and context features such that the base voxel grid consists of  $M + N - 1$  bins.

From this base voxel grid, we first extract the context voxel grid  $V^{CT}$ , visualized in yellow in Figure 3. The next step is the extraction of  $J \leq N$  correlation voxel grids  $V_j^{CR}$ , visualized in green in Figure 3, in a sliding-window fashion from the base voxel grid. While it is possible to extract up to  $J = N$  correlation voxel grids, we find in our ablation study in Section V-D.4 that the choice of  $J$  can be used to trade off performance and compute. Each of these correlation voxel grids contain information about the contrast differences in the scene at slightly different times. They will later be used to guide the search for pixel trajectories via lookup operations. We refer to the correlation voxel grids also as *views* because they each represent a distinct timestamp. If we choose to extract at least two correlation voxel grids apart from the reference timestamp  $t_r$  (i.e.  $J \geq 3$ ), we categorize the method as *multi-view*.

2) *Optional Frame-based Input*: As we show in the experimental results, the performance of the proposed method improves if image data is used in addition to events only. To do so, we acquire a reference frame  $I_r$  at  $t_r$  and a target

frame  $I_t$  at  $t_t$ . These frames contain richer information about texture at the boundary timestamps of the regressed pixel trajectories and thus simplify the correspondence search.

### C. Feature Extraction

Features are extracted from input images and voxel grids with a convolutional network. The basic architecture of the encoders follows prior work [9], [4]. The encoders extract  $D = 256$  dimensional features from their input at 1/8th of the original resolution using residual blocks with striding for downsampling the feature maps. From now on, we use  $H' = H/8$  and  $W' = W/8$  to refer to the downsampled resolution.

The context feature encoder  $f_{CT} : \mathbb{R}^{(N+3) \times H \times W} \mapsto \mathbb{R}^{D \times H' \times W'}$  concatenates the voxel grid  $V^{CT}$  with the frame  $I_r$ , if available, to extract combined features. The reference frame  $I_r$  informs the network about the absolute intensity of the reference pixels while the context voxel grid  $V^{CT}$  provides rich information about motion during the time duration of the pixel trajectories. Note that the context network only extracts features from  $I_r$  and not  $I_t$  because we require the context features to be aligned with the reference timestamp. Finally, contrast information is not explicitly considered yet, which will be provided by the correlation feature encoders.

The correlation feature encoder  $f_{CR}^V : \mathbb{R}^{M \times H \times W} \mapsto \mathbb{R}^{D \times H' \times W'}$  for event representations computes features from the  $N$  correlation voxel grids in parallel, by sharing the weights. The result of this operation are  $N$  feature maps, each assigned to the timestamp associated with the last bin of each correlation voxel grid. If in addition, a pair of frames is available, we extract image features with an additional encoder  $f_{CR}^I : \mathbb{R}^{3 \times H \times W} \mapsto \mathbb{R}^{D \times H' \times W'}$ .

### D. Multi-View Correlation Volumes

We compute correlation volumes to guide correspondence search in space and time by associating subsequent views with the reference view. This association is visualized in Figure 4 and 2. We use the features created from voxel grid  $V_0^{CR}$  and optionally  $I_r$  to create the feature maps for the reference view. For each subsequent view  $j$ , we compute features from voxel grid  $V_j^{CR}$  to compute the correlation volume as in equations (5) and (3). For the final/target view, we optionally compute a correlation volume from the boundary image features as in equations (5) and (4). We perform the computation of all correlation volumes in parallel with batched matrix multiplication.

$$C(f_{CR}^V(V_0^{CR}), f_{CR}^V(V_j^{CR})) \in \mathbb{R}^{H' \times W' \times H' \times W'}, \quad (3)$$

$$C(f_{CR}^I(I_r), f_{CR}^I(I_t)) \in \mathbb{R}^{H' \times W' \times H' \times W'}, \quad (4)$$

where

$$C_{ijkl}(a, b) = \frac{1}{\sqrt{D}} \sum_d a_{dij} \cdot b_{dkl} \quad (5)$$

### E. Iterative Multi-View Flow Updates

This section describes the update scheme that uses correlation volumes introduced in the previous section III-D.

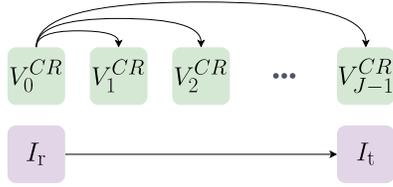


Fig. 4. Correlation-volume pairs with indicated lookup direction. The upper row shows lookups from voxel grid features and the lower row shows lookups from image features.

1) *Continuous-Time Flow*: We seek a representation of the function  $\mathbf{B}(\tau, x, y)$ , as introduced in equation (1), that generalizes the conventional displacement prediction of two-frame optical flow methods. To achieve this, we choose to use Bézier curves defined as

$$\mathbf{B}(\tau, x, y) = \sum_{i=0}^n \binom{n}{i} (1-\tau)^{n-i} \tau^i \mathbf{P}_i(x, y). \quad (6)$$

$\mathbf{B}(\tau, x, y)$  describes the displacement of pixel  $(x, y)$  for the normalized time  $\tau \in [0, 1]$ . As an example, the Bézier curves of degree  $n = 1$  simply represent a linear trajectories in space and time that is fully described by  $\mathbf{P}_1$ . The task of our method, however, is to regress the parameter set

$$\mathcal{P} = \{\mathbf{P}_1, \dots, \mathbf{P}_n\} \quad (7)$$

$\mathbf{P}_0 = \mathbf{0}$  does not have to be estimated because it defines the starting point of the trajectory that coincides with the reference pixels. The advantages of working with Bézier curves is that they are fast to evaluate and can be concisely summarized by a set of parameters  $\mathcal{P}$ .

The degree  $n$  of the Bézier curves is a fixed parameter. Our ablation study in Section V-D.2 suggests to set  $n$  equal to the number of supervision points along the trajectory for an accurate trajectory prediction.

**Iterative Bézier Updates:** The Bézier curves are initialized with  $\mathcal{P} = \{\mathbf{0}, \dots, \mathbf{0}\}$  such that  $\mathbf{B} = \mathbf{0}$ . The update block of the network estimates increments  $\Delta\mathbf{P}$  that are added to the current estimate of the Bézier curve parameters. More specifically, in each iteration from  $k$  to  $k+1$  the network updates the parameter set  $\mathcal{P}$  the following way:

$$\mathbf{P}_i^{k+1} = \mathbf{P}_i^k + \Delta\mathbf{P}_i^k, \forall i \in \{1, \dots, n\}$$

2) *Multi-View Correlation Lookup*: Similar to the original RAFT implementation, we use lookup operations to extract features from the correlation volumes. In contrast, however, we extract features from multiple correlation volumes, each associated to a unique normalized timestamp. Given the current estimate of the Bézier control points  $\mathcal{P}$ , we map each pixel  $\mathbf{x} = (x, y)$  at time  $\tau = 0$  to the estimated corresponding pixel location  $\mathbf{x}'(\tau) = (x'(\tau), y'(\tau))$  at time  $\tau$ :

$$\mathbf{x}'(\tau) = \mathbf{x} + \mathbf{B}(\tau, \mathbf{x}) \quad (8)$$

Similar to RAFT, the lookup is performed in a local neighborhood  $\mathcal{N}$  around the corresponding pixel location  $\mathbf{x}'(\tau)$ :

$$\mathcal{N}(\mathbf{x}'(\tau)) = \{\mathbf{x}'(\tau) + \mathbf{dx} \mid \mathbf{dx} \in \mathbb{Z}^2, \|\mathbf{dx}\|_\infty \leq r\} \quad (9)$$

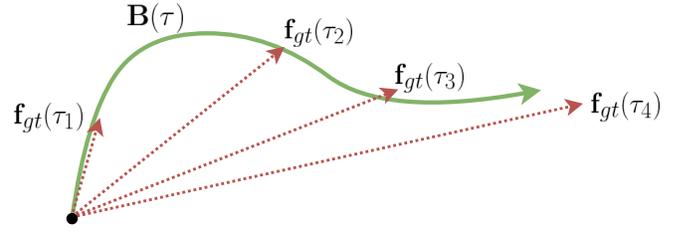


Fig. 5. Illustrative example of a Bézier prediction and  $N_k = 4$  ground truth flow maps from which the loss function (10) is computed.

Lookups are performed on all available correlation volumes with bilinear sampling. We use a constant lookup radius of  $r = 4$ , as in the original RAFT implementation, to increase the effective lookup radius. Finally, the values from the union of lookup operations are concatenated into a single feature map.

3) *Upsampling of Bézier Curves*: The Bézier control points are estimated at  $1/8$ -th of the original resolution. Since  $\mathbf{B}(\cdot)$  is linear in the control points  $\mathcal{P}$ , we can upsample the Bézier curves to the full resolution using convex upsampling [9]. As a result, the Bézier curves at the full resolution will be a learned convex combination of  $3 \times 3$  grids of Bézier curves at the lower resolution.

#### F. Supervision

We supervise the model with  $N_k$  ground truth flow maps along the trajectory, visualized in Figure 5.

$$\mathcal{L} = \frac{1}{N_k} \sum_{i=1}^{N_i} \gamma^{N_i-i} \sum_{k=1}^{N_k} \|\mathbf{f}_{gt}(\tau_k) - \mathbf{B}_i(\tau_k)\|_1 \quad (10)$$

where  $\mathbf{B}_i$  is the Bézier curve at iteration  $i$ ,  $\tau_k \in [0, 1]$  are the evaluation timestamps of the Bézier curves, and  $\gamma = 0.8$ . For a single displacement map, such as in two-frame optical flow, the corresponding parameters are  $N_k = 1$  and  $\tau_1 = 1$ .

## IV. MULTIFLOW DATASET

To the best of our knowledge, there are no publicly available datasets with dense ground truth for pixel trajectories, which is more general than pixel displacements, in combination with events and images. Therefore, we create a new synthetic dataset to evaluate the proposed approach.

### A. Data Generation Procedure

We generate 10100 training and 2000 test sequences in simulation for MultiFlow. The background images are sampled from the Flickr30K dataset [48]. The foreground objects are extracted from randomly sampled PNG images, by masking out transparent regions using their alpha channels.

The duration of each sequence is defined to be one second. We randomly sample 3 or 4 control points over similarity transformations in 2D for each foreground and background object. Interpolation with piece-wise cubic spline polynomials provide the continuous-time transformation of all objects throughout the duration of each sequence. The

	Input	EPE	AE	1PE	2PE	3PE
EV-FlowNet [3]	E	2.32	-	55.4	29.8	18.6
E-RAFT [4]	E	0.79	2.85	12.74	4.74	2.68
RAFT [9]	I	0.78	<u>2.44</u>	12.40	4.60	2.61
RAFT + GMA [1]	I	0.94	2.66	12.98	5.08	2.96
<b>Ours</b>	E	0.75	2.68	<u>11.90</u>	<u>4.41</u>	<u>2.44</u>
<b>Ours</b>	E+I	<b>0.69</b>	<b>2.42</b>	<b>9.70</b>	<b>3.42</b>	<b>1.88</b>

TABLE I

RESULTS ON DSEC-FLOW. OUR METHOD OUTPERFORMS PRIOR WORK EVEN WHEN ONLY USING EVENTS. ADDING IMAGE DATA FURTHER IMPROVES PERFORMANCE.

generated similarity transformation trajectories are used to both render images and compute ground truth pixel trajectories. To generate events, we apply the generative model of events [49] on frames rendered at 1000 frames per second on the full sequence.

**Pixel Trajectory Ground Truth:** The ground truth for the pixel trajectories is computed with respect to a reference time at 0.4 seconds. The pixel trajectory is expressed as pixel displacements at different target timestamps with respect to the original pixel location at the reference time. Overall, the dataset consists of pixel trajectory ground truth at intervals of 10 milliseconds from 0.4 seconds ( $t_r$ ) up to time 0.9 seconds ( $t_t$ ). Events outside the time interval where ground truth is available ( $0.4 \leq t \leq 0.9$ ) can be used as additional context. We refer to the supplementary material for dataset examples and a more detailed description.

## V. EXPERIMENTS

This section has four distinct purposes. First, to show quantitative real-world performance on DSEC-Flow on traditional two-view metrics. Second, to show quantitative performance on the *main task* on MultiFlow: predicting continuous-time pixel trajectories. Third, to show qualitative performance for sim2real transfer using models trained on our MultiFlow dataset. Finally, an ablation study highlights the importance of the proposed components.

**Implementation Details:** Our models are implemented in Pytorch and trained from scratch with random weights on each dataset. We use AdamW [50] with gradient clipping in the range of  $[-1, 1]$  and a OneCycle learning rate with a batch size of 3. On both datasets, we perform random horizontal and vertical flipping as well as random cropping of the input data.

### A. DSEC-Flow

DSEC-Flow [4], [51] is a driving dataset with stereo event and global shutter cameras. The purpose of this experiment is to show the quantitative performance of our approach on real-world data. However, we can only test on traditional two-view optical flow metrics because no pixel trajectory ground truth is available on DSEC-Flow. In other words, for displacement prediction, the quality of the intermediate continuous trajectory is irrelevant. Instead, only the final displacement prediction is evaluated using end-point-error

(EPE) and angular error (AE) metrics [52]. We also report the  $X$ -point error, a metric that reports the percentage of pixels with EPE higher than  $X$  pixels.

For our experiments that combine image and event data, we warp images to the event camera. Due to the small baseline between both cameras, we can warp the image to a plane at an infinite depth and re-project it into the event camera coordinate frame. The disparity is negligible.

The performance of EV-FlowNet [3] and E-RAFT is taken from Gehrig et al. [4]. We additionally train the frame-based RAFT model [9] and GMA [1], an addition over RAFT, to also compare against purely frame-based approaches.

**Implementation Details:** Our model is supervised with the 2-view version of our loss proposed in equation (10). This loss function is used because DSEC-Flow only provides ground truth for pixel displacements. We use  $M = N = 5$  bins for both the context and correlation voxel grid (see Section III-B.1) and  $J = 5$  views to compute the correlation volumes. We also experimented with a higher number of bins but did not observe significant improvements. Finally, we choose a degree of 2 for the Bézier curves mostly to account for non-linear motion due to change in depth. We train our models for 250k iterations on a single Titan RTX which takes up to 40 hours.

1) *Quantitative Evaluation:* Table I summarizes our results on the DSEC-Flow test set. Our approach, using both event data and frames, achieves 0.69 EPE which is 11.5% lower than RAFT [9] and 8 % lower than our own method using only event data. Furthermore, our purely event-based approach achieves an EPE of 0.75 which is 5% lower than the EPE of 0.79 that E-RAFT achieves. Overall the performance of E-RAFT [4] is comparable to RAFT while EV-FlowNet [3] is not competitive. In our experiments, the GMA version of RAFT does not outperform the RAFT baseline on this dataset.

These results indicate that the proposed method, even though designed to work for regressing pixel trajectories, is competitive with two-view approaches on the task of pixel displacement prediction. Note that we have not used any additional ground truth information compared to the baselines.

### B. MultiFlow

The experiments on MultiFlow assess the pixel trajectory regression capabilities. To achieve this, we introduce an extension of EPE and AE to trajectories.

$$\text{TEPE} = \frac{1}{N_k} \sum_k^{N_k} \text{EPE}(\mathbf{f}_{pred}(t_k), \mathbf{f}_{gt}(t_k)), \quad (11)$$

where  $N_k \geq 2$ . We analogously define TAE..

We train three previously published baselines for an extensive comparison. First, RAFT [9] and RAFT+GMA [1] for a comparison against frame-based approaches. Second, E-RAFT [4] for a comparison against a recent event-based approach. We focus on these architectures because they are related to our approach and achieve competitive performance on public benchmarks [53], [51].

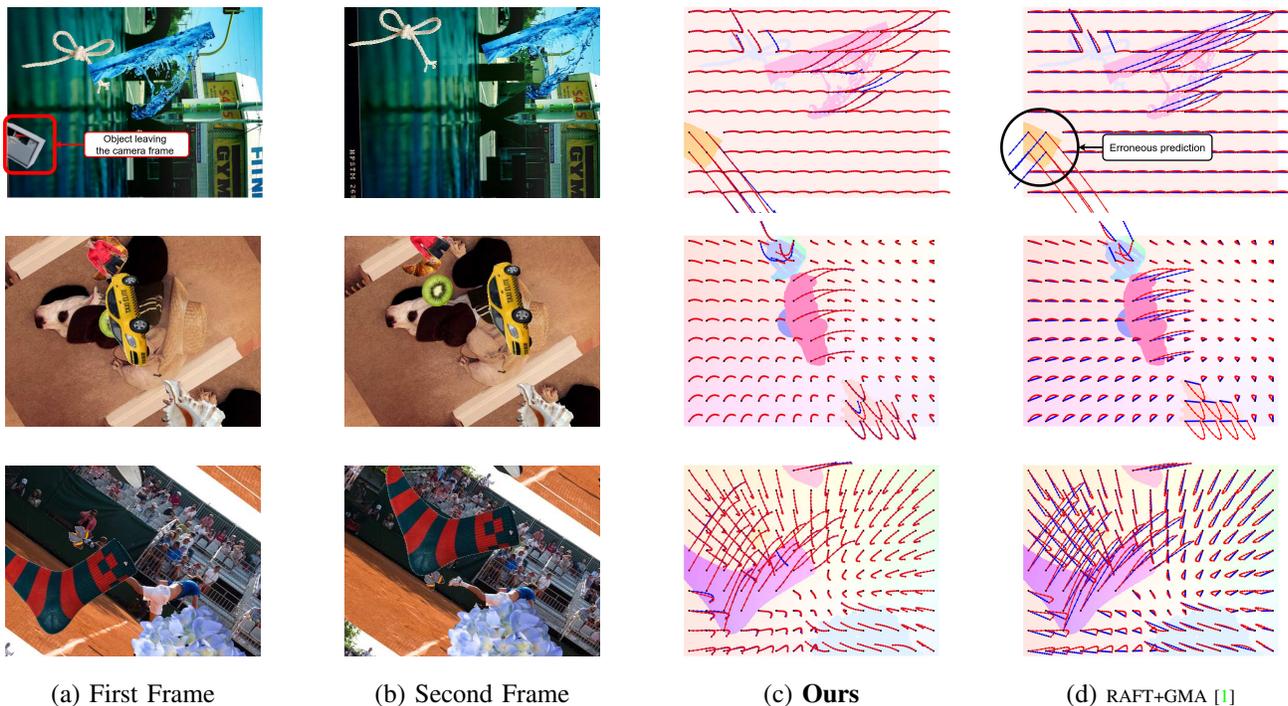


Fig. 6. Predictions of our method (c) and the strongest baseline, RAFT+GMA [1] (d), on the MultiFlow dataset. Predictions are shown in blue and the ground truth trajectory is visualized in red. The background is a colorization of the ground truth flow to highlight moving objects. Events are not shown for conciseness but are also used by our method. Best viewed in PDF form.

	Input	Trajectory Metrics		2-View Metrics	
		TEPE	TAE	EPE	AE
RAFT [9]	I	(6.89)	(19.31)	7.42	6.71
RAFT + GMA [1]	I	(5.14)	(16.35)	<b>1.47</b>	<b>1.56</b>
E-RAFT [4]	E	(6.70)	(18.44)	7.56	6.19
E-RAFT + Bézier	E	2.62	5.92	4.54	6.06
<b>Ours</b>	E	<b>1.85</b>	<b>4.61</b>	3.37	4.80
<b>Ours</b>	E+I	<b>1.29</b>	<b>3.35</b>	<b>2.27</b>	<b>3.19</b>

TABLE II

RESULTS ON MULTIFLOW. TEPE AND TAE ARE THE TRAJECTORY VERSIONS OF EPE AND AE. METRICS IN BRACKETS ARE COMPUTED USING A LINEAR MOTION MODEL, HIGHLIGHTING THAT THESE METHODS ARE NOT ORIGINALLY DESIGNED FOR ACCURATE TRAJECTORY PREDICTION. E-RAFT + BÉZIER REPRESENTS OUR BASELINE MODEL CONSISTING OF E-RAFT THAT ESTIMATES BÉZIER CURVES INSTEAD OF PIXEL DISPLACEMENTS.

**Implementation Details:** We train our method on the loss defined by equation (10) while the two-view approaches are trained on the two-view version of the loss [9]<sup>1</sup>. We supervise our methods with 10 flow maps along the trajectory using loss (10) and set the Bézier curve degree to 10 according to section III-E.1. The ablation study in Section V-D.2 clarifies the relationship between the loss function (10) and the degree of the Bézier curve.

<sup>1</sup>It is also possible to use equation (10) to supervise two-view methods, but it results in an unfavorable trade-off between trajectory end-point-error (TEPE) and end-point-error (EPE).

For the following experiments, we build  $J = 6$  correlation voxel grids representing the views at regular time intervals starting from 0.4 seconds up to 0.9 seconds (i.e. the duration where we have to predict the pixel trajectories in the dataset). We discretize the time into regular bins resulting in  $N = 41$  bins for the context voxel grid and  $M = 25$  bins for the correlation voxel grids. Later in section V-D.5, we show that a coarse discretization with fewer bins leads to suboptimal performance. We observed that the reason for this is the overwriting of polarities in the voxel grid which we prevent by choosing a more fine-grained temporal discretization.

We train our models for 200k iterations on a single Titan RTX which takes up to 80 hours.

1) *Quantitative Evaluation:* In addition to E-RAFT and RAFT variants using frames, we also compare against a baseline that consists of the E-RAFT approach that estimates Bézier curves instead of pixel displacements, as in the original work. This approach, named *E-RAFT + Bézier*, is better suited for non-linear flow estimation than prior work. Table II summarizes the quantitative results. Trajectory metrics in brackets indicate approaches that predict pixel displacements (RAFT, RAFT + GMA, and E-RAFT), for which we employ linear interpolation to compute the trajectory metrics. As expected, the results suggest that these approaches are unsuitable for accurate pixel trajectory estimation. The *E-RAFT + Bézier* baseline clearly improves not only the trajectory metrics but also the 2-view metrics compared to these baselines. Our proposed model substantially reduces the TEPE further from 2.62 to 1.85 by using correlation features to guide the trajectory estimation. Finally, our approach benefits

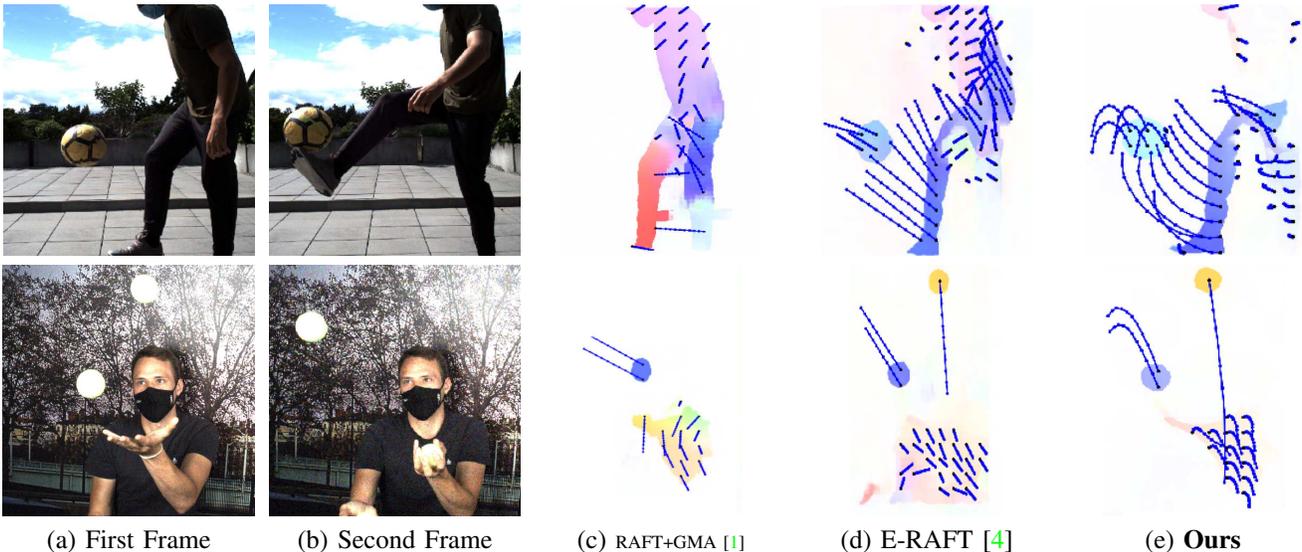


Fig. 7. Predictions of our method using only event data (e), E-RAFT [4] using events (d), and RAFT+GMA [1] using frames (c), on the HS-ERGB [6] dataset. Predictions are shown in blue. The background is a colorization flow to highlight moving objects. Events are not shown for conciseness.

from combining events and frames: TEPE is reduced from 1.85 to 1.29 by 30% while the EPE also decreases by 32%.

When purely comparing 2-view metrics, RAFT + GMA [1] achieves the lowest error by extending RAFT with a self-attention mechanism on the context features to aggregate motion features. However, RAFT+GMA is not designed for trajectory prediction and, therefore, falls short in TEPE and TAE compared to our approach.

2) *Qualitative Analysis*: Figure 6 illustrates the predictions of our model, using events and frames, in comparison to RAFT+GMA [1]. Our approaches successfully predict the continuous trajectory, even if the objects leave the field of view. This is not the case for the two-view method RAFT+GMA, as seen in the first row of Figure 6. Evidently, it fails at predicting the motion of objects that are missing in the second view because it is impossible to establish a match. The second row shows that our method can accurately predict pixel trajectories while RAFT+GMA does well in predicting the final pixel displacement.

### C. Qualitative Sim2Real Results

To qualitatively assess the real-world capabilities of our method, we use our purely event-based model, trained only on the simulation dataset MultiFlow, and show predictions on the HS-ERGB and BS-ERGB datasets [6], [7] in Figure 1 and 7. Although our method is only trained on MultiFlow, it can correctly predict non-linear motion of moving objects while the frame-based baseline [1] fails due to large motion and ambiguities in the input frames.

### D. Ablation Study

This section examines the influence of the main contributions of our model as well as additional design choices that contribute to the performance of the method. We perform the ablation studies on the MultiFlow because we require accurate pixel trajectory ground truth that is not available on DSEC or other event-based datasets.

Input	Trajectory Metrics		2-View Metrics		
	TEPE	TAE	EPE	AE	
w/o Bézier	E+I	6.36	18.37	6.66	5.18
w/o multi-view	E+I	1.81	4.25	2.68	3.51
<u>Reference</u>	E+I	<u>1.29</u>	<u>3.35</u>	<u>2.27</u>	<u>3.19</u>
w/o Bézier	E	7.31	20.07	9.10	8.85
w/o multi-view	E	2.62	5.92	4.54	6.06
<u>Reference</u>	E	<u>1.85</u>	<u>4.61</u>	<u>3.37</u>	<u>4.80</u>

TABLE III

ABLATION EXPERIMENTS ON MULTIFLOW. FINAL MODEL SETTINGS ARE UNDERLINED.

1) *Main Components*: One of the main contributions of this paper are the iterative Bézier curve regression and the correlation lookup applied at multiple timesteps/views along the trajectories. We ablate both components using the full model described in Section III as reference and summarize the results in Table III.

First, we remove the Bézier curve regression and note a drastic drop in performance on all metrics. Removing the Bézier curves is equivalent to assuming linear motion. Linear motion is not representative for this dataset and impairs the utility of the correlation lookup. For example, if the model is constrained to predict a linear trajectory during non-linear motion, the correlation lookup will not be accurately placed along the trajectory. This leads to correlation features that are not informative.

Second, we remove the correlation lookups between the reference view and intermediate views, denoted as the w/o multi-view baseline in Table III, and note that both TEPE and EPE increase between 18% and 42% with respect to the event-based and hybrid reference model. We conclude that both components are essential for achieving the best results.

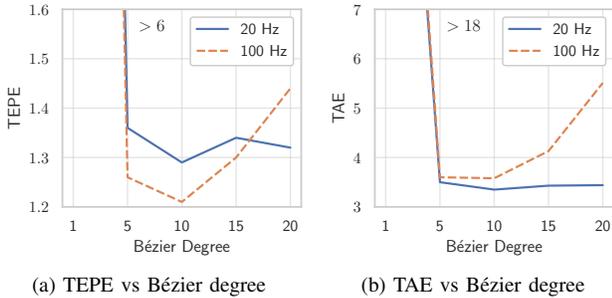


Fig. 8. **Trajectory metrics vs Bézier curve degree.** (a) shows the trajectory EPE and (b) the trajectory AE as a function of the Bézier degree. We evaluate both at 20 Hz (every 50 milliseconds) and at 100 Hz (every 10 ms). In both cases we can observe overfitting of the predicted trajectories when we evaluate at 100 Hz. This can be explained with the fact that the loss function was only applied every 50 milliseconds during training.

2) *Degree of the Bézier Curves:* The goal of this experiment is to answer the question: *Which Bézier curve degree is appropriate for our model?*

To answer this question, we train 5 different models with Bézier curve degrees 1, 5, 10, 15, and 20 using frames and events as input to the model. The supervision on pixel trajectories at training time is at 20 Hz, that we supervise the model on the pixel trajectory every 50 milliseconds.

Intuitively, one could surmise that a higher Bézier degree always improves the performance. Figure 8 indeed shows that a higher Bézier degree improves performance, when the models are evaluated at 20 Hz. The models were also trained with a supervision signal at 20 Hz, that is with supervision at regular time intervals of 50 milliseconds. Interestingly, an evaluation of these models at a higher frequency of 100 Hz reveals a decrease in performance for the model with a Bézier degree higher than 10. The reason for that is that we observe a test-time overfitting of the predicted trajectory to the frequency with which the model was supervised. We do not observe this overfitting issue at 20 Hz evaluation frequency because the model is able to accurately predict the trajectory at this frequency. This observation is reminiscent of the overfitting phenomenon in the much simpler polynomial regression problem in statistics.

The answer to the original question is: *For accurate pixel trajectory predictions, the Bézier degree should be chosen equal to the number of regular supervision points along the pixel trajectories.* The degree can be set higher than the number of supervision points to potentially improve performance on the timestamps of interest. However, the pixel trajectories may not be as accurate anymore. An example for this is our DSEC-Flow experiments where we have only 1 supervision point (2-view) but found that a Bézier degree of 2 improves the performance on the 2-view metric. This is a reasonable approach as long as the downstream application does not require accurate pixel trajectories.

3) *Loss Function:* This experiment examines the influence of the trajectory loss on trajectory and 2-view metrics. We train two additional models with the 2-view loss; one using only events and a second one using both frames and events. Table IV summarizes the results.

For both input variations, the error metrics substantially decrease when the model is trained with the trajectory loss. This reduction is evident for both the trajectory metrics and also the 2-view metrics. This indicates, that the trajectory loss function is a better choice for the proposed method than the 2-view loss. Finally, the results in Table IV indicate that we could further reduce the errors of the DSEC-Flow experiment with the appropriate ground truth.

	Input	Trajectory Metrics		2-View Metrics	
		TEPE	TAE	EPE	AE
2-View Loss	E+I	16.47	20.57	3.72	5.48
<u>Trajectory Loss</u>	<u>E+I</u>	<u>1.29</u>	<u>3.35</u>	<u>2.27</u>	<u>3.19</u>
2-View Loss	E	16.99	21.83	5.95	8.57
<u>Trajectory Loss</u>	<u>E</u>	<u>1.85</u>	<u>4.61</u>	<u>3.37</u>	<u>4.80</u>

TABLE IV

THE TRAJECTORY LOSS REFERS TO THE LOSS FUNCTION OF EQUATION (10) OF THE MAIN PAPER, WITH  $N_k > 1$ . THE 2-VIEW LOSS REFERS TO THE SAME LOSS FUNCTION WITH  $N_k = 1$ . TRAINING OUR METHOD WITH THE TRAJECTORY LOSS LEADS TO DRASTICALLY LOWER ERRORS EVEN FOR THE 2-VIEW METRICS.

4) *Number of Correlation Lookups:* Table V shows that increasing the number of correlation lookups improves performance consistently on all metrics. However, substantially increasing the number of correlation volumes incurs higher memory consumption and increases computational demand. This result suggests a trade-off between performance and compute and memory requirements. Note that the number of correlation volumes is the number of correlation lookups (in time) + 1 because the correlation volume at the reference time has to be taken into account as well.

# Correlation Lookups	Trajectory Metrics		2-View Metrics	
	TEPE	TAE	EPE	AE
1	1.81	4.25	2.68	3.51
3	1.35	3.43	2.34	3.21
<u>5</u>	<u>1.29</u>	<u>3.35</u>	<u>2.27</u>	<u>3.19</u>

TABLE V

INCREASING THE NUMBER OF CORRELATION VOLUMES/LOOKUPS IN TIME IMPROVES PERFORMANCE. THE UNDERLINED ROW REFERS TO THE REFERENCE MODEL USED IN THE EXPERIMENTS.

5) *Voxel Grid Temporal Resolution:* The voxel grid serves as a discrete representation of event data, inevitably leading to some loss of information during its formation. In particular, a decrease in the number of temporal bins not only reduces the number of input channels but also intensifies the quantization of the signal. In our study, detailed in Table V-D.4, we investigate how the temporal resolution of the voxel grid affects performance on the MultiFlow dataset. The bin size in this context refers to the temporal interval between consecutive bins. Consistent with expectations, Table V-D.4

demonstrates that a finer bin size, corresponding to higher temporal resolution, results in improved performance. Note that a bin size of 0.0125 represents the voxel grid resolution, as detailed in section V-B, which was used for the main results on the MultiFlow dataset.

Bin Size	Trajectory Metrics		2-View Metrics	
	TEPE	TAE	EPE	AE
0.1	1.81	4.53	2.87	3.97
0.05	1.51	3.82	2.61	3.63
0.025	1.40	3.56	2.45	3.41
<u>0.0125</u>	<u>1.29</u>	<u>3.35</u>	<u>2.27</u>	<u>3.19</u>

TABLE VI

DECREASING THE BIN SIZE OF THE VOXEL GRID IMPROVES PERFORMANCE. EACH BIN IS ASSIGNED A SPECIFIC TIMESTAMP AND THE BIN SIZE REFERS TO THE  $\Delta t$  BETWEEN THE BINS. THE UNDERLINED ROW REFERS TO THE REFERENCE MODEL USED IN THE EXPERIMENTS.

### E. Runtime Analysis

Table VII shows inference time in milliseconds, parameter count in millions of parameters and memory consumption at inference time on samples from the MultiFlow dataset. Overall, our proposed method has higher demands on the presented metrics compared to RAFT [9], which is expected because it is a generalization of the RAFT architecture.

	Input	Inference time [ms]	Params	Memory [GB]
RAFT [9]	I	52	5.3 M	1.20
E-RAFT [4]	E	56	5.3 M	1.20
Ours	E	77	5.6 M	1.65
Ours	E+I	84	5.9 M	1.65

TABLE VII

COMPARISON OF INFERENCE TIME, PARAMETER COUNT AND MEMORY CONSUMPTION ON MULTIFLOW. THESE NUMBERS HAVE BEEN OBTAINED ON A TITAN RTX GPU WITH AN IMPLEMENTATION USING PYTORCH VERSION 1.12.1

## VI. LIMITATIONS

The current implementation computes the dot products between all paired features at once, which incurs  $\mathcal{O}(N^2)$  space complexity where N is the number of spatial features. This limitation can be overcome by computing dot products on-demand, as originally proposed by Teed et al. [9]. Another limitation of the proposed method is that we cannot compute the pixel trajectories for a sequence longer than the current input to the model. We can approximate longer trajectories by concatenating multiple shorter sequences via interpolation at the reference timestamps. A more elegant solution could be achieved by designing a recurrent model that can track pixels without assuming a fixed sequence length or duration. Finally, the number of correlation lookups is a fixed parameter that cannot be changed at inference time. As a result, we must define at training time how many correlation lookups can be afforded during inference time.

## VII. CONCLUSION

We have introduced a method for estimating continuous-time pixel trajectories from events and frames. The proposed approach generalizes the frame-based [9] and event-based [4] RAFT architecture by regressing Bézier curves and uses the pixel trajectories to extract correlation features along the temporal axis. Our experimental results demonstrate that the proposed method can accurately predict continuous pixel trajectories while at the same time outperforming strong baselines not only in simulation but also on real data. Finally, our sim2real results suggest that the new MultiFlow dataset can also be used to pretrain pixel trajectory regression models for downstream applications.

## ACKNOWLEDGMENT

This work was supported by the National Centre of Competence in Research (NCCR) Robotics (grant agreement No. 51NF40-185543) through the Swiss National Science Foundation (SNSF), and the European Research Council (ERC) under grant agreement No. 864042 (AGILEFLIGHT).

## REFERENCES

- [1] S. Jiang, D. Campbell, Y. Lu, H. Li, and R. Hartley, "Learning to estimate hidden motions with global motion aggregation," in *Int. Conf. Comput. Vis. (ICCV)*, 2021.
- [2] F. Guney, L. Sevilla-Lara, D. Sun, and J. Wulff, "What is optical flow for?": Workshop results and summary," in *Eur. Conf. Comput. Vis. Workshops (ECCVW)*, 2018.
- [3] A. Z. Zhu, L. Yuan, K. Chaney, and K. Daniilidis, "EV-FlowNet: Self-supervised optical flow estimation for event-based cameras," in *Robotics: Science and Systems (RSS)*, 2018.
- [4] M. Gehrig, M. Millhäusler, D. Gehrig, and D. Scaramuzza, "E-RAFT: Dense optical flow from event cameras," in *3D Vision (3DV)*, 2021.
- [5] G. Gallego, T. Delbruck, G. Orchard, C. Bartolozzi, B. Taba, A. Censi, S. Leutenegger, A. Davison, J. Conradt, K. Daniilidis, and D. Scaramuzza, "Event-based vision: A survey," *IEEE Trans. Pattern Anal. Mach. Intell.*, 2020.
- [6] S. Tulyakov, D. Gehrig, S. Georgoulis, J. Erbach, M. Gehrig, Y. Li, and D. Scaramuzza, "TimeLens: Event-based video frame interpolation," *IEEE Conf. Comput. Vis. Pattern Recog. (CVPR)*, 2021.
- [7] S. Tulyakov, A. Bochicchio, D. Gehrig, S. Georgoulis, Y. Li, and D. Scaramuzza, "Time lens++: Event-based frame interpolation with parametric non-linear flow and multi-scale fusion," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2022, pp. 17 755–17 764.
- [8] A. Dosovitskiy, P. Fischer, E. Ilg, P. Hausser, C. Hazirbas, V. Golkov, P. van der Smagt, D. Cremers, and T. Brox, "FlowNet: Learning optical flow with convolutional networks," in *Int. Conf. Comput. Vis. (ICCV)*, December 2015.
- [9] Z. Teed and J. Deng, "Raft: Recurrent all-pairs field transforms for optical flow," in *Eur. Conf. Comput. Vis. (ECCV)*, 2020.
- [10] H. Xu, J. Yang, J. Cai, J. Zhang, and X. Tong, "High-resolution optical flow from 1d attention and correlation," in *Int. Conf. Comput. Vis. (ICCV)*, October 2021, pp. 10 498–10 507.
- [11] D. Sun, X. Yang, M.-Y. Liu, and J. Kautz, "PWC-Net: CNNs for optical flow using pyramid, warping, and cost volume," in *IEEE Conf. Comput. Vis. Pattern Recog. (CVPR)*, 2018.
- [12] E. Ilg, N. Mayer, T. Saikia, M. Keuper, A. Dosovitskiy, and T. Brox, "FlowNet 2.0: Evolution of optical flow estimation with deep networks," in *IEEE Conf. Comput. Vis. Pattern Recog. (CVPR)*, 2017, pp. 1647–1655.
- [13] R. Jonschkowski, A. Stone, J. T. Barron, A. Gordon, K. Konolige, and A. Angelova, "What matters in unsupervised optical flow," in *Eur. Conf. Comput. Vis. (ECCV)*, 2020.
- [14] J. Wulff and M. J. Black, "Efficient sparse-to-dense optical flow estimation using a learned basis and layers," in *IEEE Conf. Comput. Vis. Pattern Recog. (CVPR)*, June 2015.

- [15] D. Gadot and L. Wolf, "Patchbatch: A batch augmented loss for optical flow," in *IEEE Conf. Comput. Vis. Pattern Recog. (CVPR)*, June 2016.
- [16] C. Bailer, K. Varanasi, and D. Stricker, "Cnn-based patch matching for optical flow with thresholded hinge embedding loss," in *IEEE Conf. Comput. Vis. Pattern Recog. (CVPR)*, July 2017.
- [17] F. Zhang, O. J. Woodford, V. A. Prisacariu, and P. H. Torr, "Separable flow: Learning motion cost volumes for optical flow estimation," in *Int. Conf. Comput. Vis. (ICCV)*, October 2021, pp. 10 807–10 817.
- [18] A. Jaegle, S. Borgeaud, J.-B. Alayrac, C. Doersch, C. Ionescu, D. Ding, S. Koppula, D. Zoran, A. Brock, E. Shelhamer, O. Hénaff, M. M. Botvinick, A. Zisserman, O. Vinyals, and J. Carreira, "Perceiver io: A general architecture for structured inputs and outputs," 2021.
- [19] S. Zhao, Y. Sheng, Y. Dong, E. I. Chang, Y. Xu *et al.*, "Maskflownet: Asymmetric feature matching with learnable occlusion mask," in *IEEE Conf. Comput. Vis. Pattern Recog. (CVPR)*, 2020.
- [20] Y. Zheng, M. Zhang, and F. Lu, "Optical flow in the dark," in *IEEE Conf. Comput. Vis. Pattern Recog. (CVPR)*, June 2020.
- [21] A. Dosovitskiy, P. Fischer, E. Ilg, P. Häusser, C. Hazırbaş, V. Golkov, P. van der Smagt, D. Cremers, and T. Brox, "FlowNet: Learning optical flow with convolutional networks," in *Int. Conf. Comput. Vis. (ICCV)*, 2015, pp. 2758–2766.
- [22] P. Godet, A. Boulch, A. Plyer, and G. L. Besnerais, "Starflow: A spatiotemporal recurrent cell for lightweight multi-frame optical flow estimation," in *IEEE Int. Conf. Pattern Recog. (ICPR)*, 2021, pp. 2462–2469.
- [23] J. Janai, F. Güney, A. Ranjan, M. Black, and A. Geiger, "Unsupervised learning of multi-frame optical flow with occlusions," in *Eur. Conf. Comput. Vis. (ECCV)*, September 2018.
- [24] D. Maurer and A. Bruhn, "Proflow: Learning to predict optical flow," in *British Mach. Vis. Conf. (BMVC)*, 06 2018.
- [25] J. Janai, F. Güney, J. Wulff, M. Black, and A. Geiger, "Slow flow: Exploiting high-speed cameras for accurate and diverse optical flow reference data," in *IEEE Conf. Comput. Vis. Pattern Recog. (CVPR)*, 2017.
- [26] S. Ricco and C. Tomasi, "Dense lagrangian motion estimation with occlusions," in *IEEE Conf. Comput. Vis. Pattern Recog. (CVPR)*, 2012, pp. 1800–1807.
- [27] R. Garg, A. Roussos, and L. Agapito, "A variational approach to video registration with subspace constraints," *Int. J. Comput. Vis.*, vol. 104, no. 3, pp. 286–314, 2013.
- [28] Z. Ren, O. Gallo, D. Sun, M.-H. Yang, E. B. Sudderth, and J. Kautz, "A fusion approach for multi-frame optical flow estimation," in *IEEE Winter Conf. Appl. Comput. Vis. (WACV)*, 2019, pp. 2077–2086.
- [29] A. W. Harley, Z. Fang, and K. Fragkiadaki, "Particle video revisited: Tracking through occlusions using point trajectories," in *ECCV*, 2022.
- [30] R. Benosman, S.-H. Ieng, C. Clercq, C. Bartolozzi, and M. Srinivasan, "Asynchronous frameless event-based optical flow," *Neural Netw.*, vol. 27, pp. 32–37, 2012.
- [31] D. Gehrig, H. Rebecq, G. Gallego, and D. Scaramuzza, "EKLT: Asynchronous photometric feature tracking using events and frames," *Int. J. Comput. Vis.*, 2019.
- [32] B. D. Lucas and T. Kanade, "An iterative image registration technique with an application to stereo vision," in *Int. Joint Conf. Artificial Intell. (IJCAI)*, 1981, pp. 674–679.
- [33] D. J. Butler, J. Wulff, G. B. Stanley, and M. J. Black, "A naturalistic open source movie for optical flow evaluation," in *Eur. Conf. Comput. Vis. (ECCV)*, Oct. 2012.
- [34] E. Mueggler, C. Forster, N. Baumli, G. Gallego, and D. Scaramuzza, "Lifetime estimation of events from dynamic vision sensors," in *IEEE Int. Conf. Robot. Autom. (ICRA)*, 2015, pp. 4874–4881.
- [35] P. Bardow, A. J. Davison, and S. Leutenegger, "Simultaneous optical flow and intensity estimation from an event camera," in *IEEE Conf. Comput. Vis. Pattern Recog. (CVPR)*, 2016, pp. 884–892.
- [36] L. Pan, M. Liu, and R. Hartley, "Single image optical flow estimation with an event camera," *IEEE Conf. Comput. Vis. Pattern Recog. (CVPR)*, 2020.
- [37] A. Z. Zhu, L. Yuan, K. Chaney, and K. Daniilidis, "Unsupervised event-based learning of optical flow, depth, and egomotion," in *IEEE Conf. Comput. Vis. Pattern Recog. (CVPR)*, 2019.
- [38] C. Lee, A. K. Kosta, A. Z. Zhu, K. Chaney, K. Daniilidis, and K. Roy, "Spike-flownet: event-based optical flow estimation with energy-efficient hybrid neural networks," in *Eur. Conf. Comput. Vis. (ECCV)*, 2020.
- [39] F. Paredes-Vallés and G. C. de Croon, "Back to event basics: Self-supervised learning of image reconstruction for event cameras via photometric constancy," in *IEEE Conf. Comput. Vis. Pattern Recog. (CVPR)*, 2021.
- [40] C. Ye, A. Mitrokhin, C. Fermüller, J. A. Yorke, and Y. Aloimonos, "Unsupervised learning of dense optical flow, depth and egomotion with event-based sensors," in *IEEE/RSJ Int. Conf. Intell. Robot. Syst. (IROS)*, 2020.
- [41] C. Kerl, J. Stückler, and D. Cremers, "Dense continuous-time tracking and mapping with rolling shutter RGB-D cameras," in *Int. Conf. Comput. Vis. (ICCV)*, 2015.
- [42] E. Mueggler, G. Gallego, H. Rebecq, and D. Scaramuzza, "Continuous-time visual-inertial odometry for event cameras," *IEEE Trans. Robot.*, vol. 34, no. 6, pp. 1425–1440, Dec. 2018.
- [43] H. Seok and J. Lim, "Robust feature tracking in DVS event stream using bézier mapping," in *IEEE Winter Conf. Appl. Comput. Vis. (WACV)*, 2020.
- [44] I. Alzugaray López and M. Chli, "Haste: Multi-hypothesis asynchronous speeded-up tracking of events," in *British Mach. Vis. Conf. (BMVC)*, 2020.
- [45] A. Z. Zhu, N. Atanasov, and K. Daniilidis, "Event-based feature tracking with probabilistic data association," in *IEEE Int. Conf. Robot. Autom. (ICRA)*, 2017.
- [46] N. Mayer, E. Ilg, P. Häusser, P. Fischer, D. Cremers, A. Dosovitskiy, and T. Brox, "A large dataset to train convolutional networks for disparity, optical flow, and scene flow estimation," in *IEEE Conf. Comput. Vis. Pattern Recog. (CVPR)*, 2016.
- [47] D. Sun, D. Vlasic, C. Herrmann, V. Jampani, M. Krainin, H. Chang, R. Zabih, W. T. Freeman, and C. Liu, "Autoflow: Learning a better training set for optical flow," in *IEEE Conf. Comput. Vis. Pattern Recog. (CVPR)*, 2021.
- [48] P. Young, A. Lai, M. Hodosh, and J. Hockenmaier, "From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions," *Transactions of the Association for Computational Linguistics*, 2014.
- [49] H. Rebecq, D. Gehrig, and D. Scaramuzza, "ESIM: an open event camera simulator," in *Conf. on Robotics Learning (CoRL)*, 2018.
- [50] I. Loshchilov and F. Hutter, "Decoupled weight decay regularization," in *Int. Conf. Learn. Representations (ICLR)*, 2019. [Online]. Available: <https://openreview.net/forum?id=Bkg6RiCqY7>
- [51] M. Gehrig, W. Aarents, D. Gehrig, and D. Scaramuzza, "Dsec: A stereo event camera dataset for driving scenarios," *IEEE Robot. Autom. Lett.*, 2021.
- [52] S. Baker, D. S. J. Lewis, S. Roth, M. J. Black, and R. Szeliski, "A database and evaluation methodology for optical flow," *Int. J. Comput. Vis.*, vol. 92, no. 1, pp. 1–31, 2011.
- [53] D. J. Butler, J. Wulff, G. B. Stanley, and M. J. Black, "A naturalistic open source movie for optical flow evaluation," in *Eur. Conf. Comput. Vis. (ECCV)*, 2012.

# Supplementary Material

## I. INFLUENCE OF TRAJECTORY LENGTH ON THE TRAJECTORY ERROR

In Figure 1, we examine how the trajectory length influences the trajectory error on the MultiFlow. As expected, the trajectory error increases when the trajectory becomes longer because the task complexity increases. For example, longer trajectories correlate with more complex motion patterns and occlusions which are particularly challenging for optical flow methods.

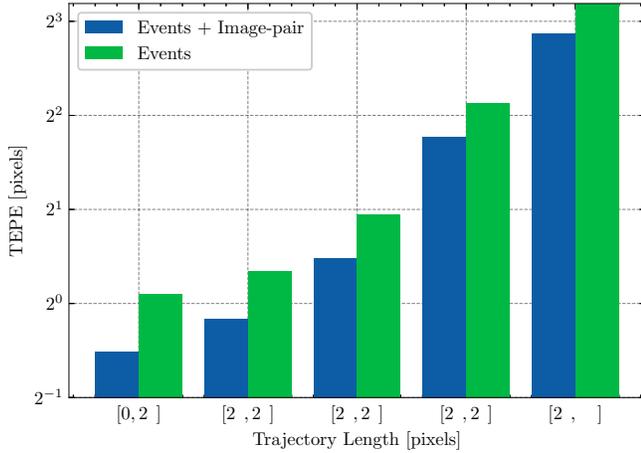


Fig. 1. Trajectory end-point-error (TEPE) as a function of the trajectory length. We use the  $\log_2$  scale to better capture the full range of trajectory length and errors. The error increases with the trajectory length, and image data generally lowers the error.

## II. NETWORK ARCHITECTURE

Figure 2 shows the neural network architecture used throughout this paper. The architecture is based on the RAFT neural network [1] with modifications highlighted in **red**. We use the same encoder architecture to extract context features from voxel grids (events) and, optionally, the image at the reference timestamp, as well as correlation features from both modalities. This encoder architecture is identical to the RAFT context and correlation encoder. Instead, the update block requires a few modifications:

- Instead of encoding pixel displacements (flow), we encode the Bézier curve parameters  $\mathbf{P}_{k-1}$  from the previous output of the update block.
- We add a context encoder to extract motion features directly from event data and optionally use the image context encoder.

- We predict  $2n$  scalars per feature of the update block to estimate the increment of  $n$  two-dimensional Bézier curve control points  $\Delta\mathbf{P}_k$ .

The upsampling module, omitted for brevity, is identical to the one used in RAFT.

## III. MULTIFLOW DATASET GENERATION

This section provides a more formal and detailed overview of the data generation process that was used to create the MultiFlow dataset.

The duration of each sequence in the dataset is fixed to one second. We sample either 3 or 4 (random choice of 50% probability each) control points in the space of pixel translation, scale, and rotation of the foreground (FG) and background (BG) objects in the image plane. The time assigned to control points is sampled from the range  $[0, 1]$  while ensuring strictly increasing values to maintain the temporal order.

The control points are sampled according to the following discrete-time process with parameters shown in Table I:

$$X_{k+1} = \begin{cases} \hat{\gamma} \cdot \hat{X}_{k+1}^{\text{Det}} + (\hat{\gamma} - 1) \cdot \hat{X}_{k+1}^{\text{Stoch}} & (\hat{\alpha} = 0) \wedge (\hat{\beta} = 0) \\ X_k & (\hat{\alpha} = 1) \vee (\hat{\beta} = 1) \end{cases} \quad (1)$$

where

$$\hat{X}_{k+1}^{\text{Det}} = X_k + \frac{t_{k+1} - t_k}{t_k - t_{k-1}} (X_k - X_{k-1}) \quad (2)$$

refers to the deterministic part of the process that uses a constant velocity model. It is inspired by the law of conservation of momentum from classical mechanics.  $t_k$  is the time associated with control point  $k$ . To introduce more variability in the data generation process, we introduce a stochastic process:

For translation and rotation:

$$\hat{X}_{k+1}^{\text{Stoch}} = X_k + \Delta X \quad \Delta X \sim \text{Uni}(-\theta, \theta) \quad (3)$$

For scale:

$$\hat{X}_{k+1}^{\text{Stoch}} = X_k \cdot \Delta X \quad \Delta X \sim (1 + \delta_0)^{(2\delta_1 - 1)}, \quad (4)$$

$$\delta_0 \sim \text{Uni}(0, \theta), \quad \delta_1 \sim \text{Bern}(0.5)$$

*Uni* stands for the continuous uniform distribution and *Bern* refers to the Bernoulli distribution.

$\hat{\alpha} \sim \text{Bern}(\alpha)$  models the probability that the similarity transformation (translation, rotation *and* scale) of the object remains constant during the whole sequence.

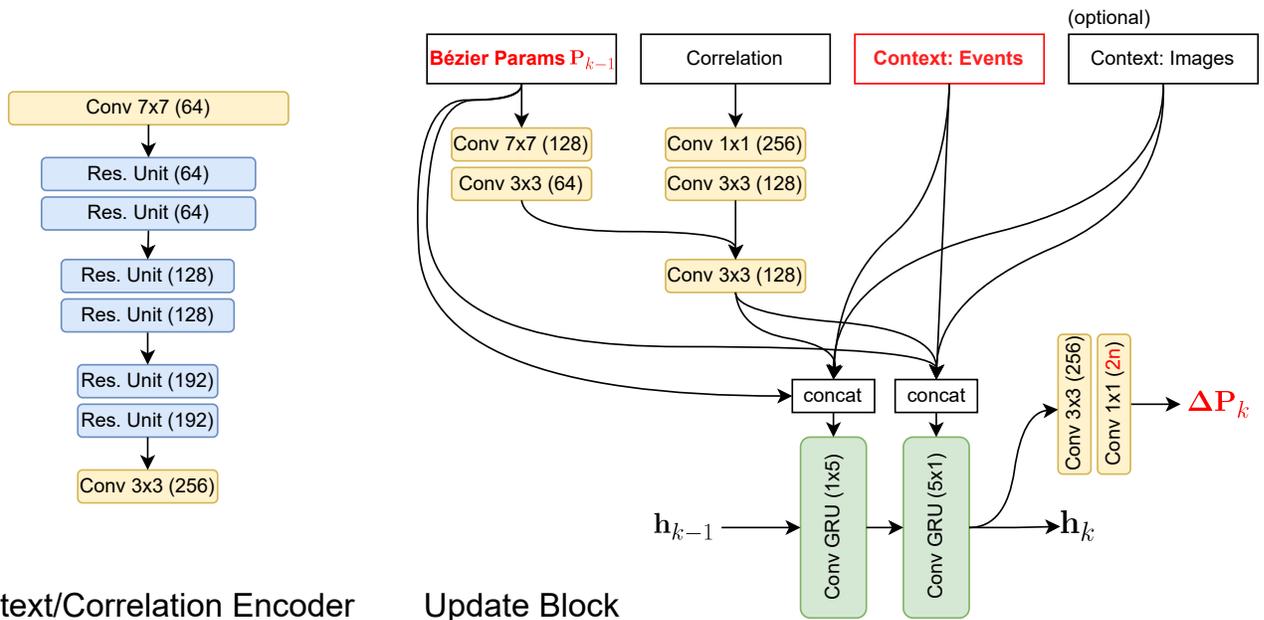


Fig. 2. Neural network architecture adapted from RAFT to predict Bézier curves. Modifications in the update block are highlighted in **red**. Context and correlation encoder architectures remain identical to RAFT, with an added context encoder for event data.  $n$  refers to the degree of the Bézier curve.

$\hat{\beta} \sim \text{Bern}(\beta)$  is sampled separately for each component (translation, rotation *or* scale) of the similarity transformation and models the probability that the component remains constant during the sequence.

$\hat{\gamma} \sim \text{Uni}(0, \gamma)$  is sampled separately for each component (translation, rotation *or* scale) of the similarity transformation and models the degree of the constant velocity component. For  $\hat{\gamma} = 1$ , the new control point is computed deterministically with a constant velocity model according to equation (2).  $\hat{\gamma} = 0$  corresponds to a purely stochastic update according to equations (3) or (4).  $\hat{\gamma} \in (0, 1)$  results in a convex combination of both as shown in equation (1).

Component	$\alpha$	$\beta$	$\gamma$	$\theta$
Translation BG	0	0.8	30	
Rotation BG	0.1	0.7	0.6	10
Scale BG		0.4	0.3	0.15
Translation FG		0	0.9	120
Rotation FG	0	0.3	0.6	30
Scale FG		0.3	0.3	0.30

TABLE I

PARAMETERS OF THE DISTRIBUTION OF THE TRANSFORMATION COEFFICIENTS FOR FOREGROUND OBJECTS (FG) AND BACKGROUND (BG). THE VALUES OF  $\theta$  ARE SHOWN IN *pixels* FOR TRANSLATION, IN *degree* FOR ROTATION AND IN *dimensionless quantity* FOR SCALE. THE VALUES OF  $\alpha$  AFFECT THE OVERALL SIMILARITY TRANSFORMATION AND THEREFORE BELONG TO TRANSLATION, ROTATION AND SCALE SIMULTANEOUSLY.

transformations via cubic spline interpolation on the  $K \in \{3, 4\}$  acquired control points  $X_0, \dots, X_{K-1}$ . It would be possible to sample a higher number of control points, but we found that 4 control points are sufficient to capture interesting non-linear motion that we could also observe in real-world applications of our algorithm (see Fig. 1 and 7 in the paper).

**Pixel Trajectory Ground Truth:** The ground truth for the pixel trajectories is computed with respect to a reference time at 0.4 seconds. The pixel trajectory is expressed as pixel displacements at different target timestamps with respect to the original pixel location at the reference time. Overall, the dataset consists of pixel trajectory ground truth at intervals of 10 milliseconds up to time 0.9 seconds. The first 0.4 and last 0.1 seconds can be used as additional context.

**Visual Examples:** Figure 3 shows examples of frames at reference timestamp in column (a) and target timestamp in column (b) as well as corresponding predictions of our model trained on the dataset. Frames between the reference and target timestamps are used to generate events that are used by our model to predict the continuous-time pixel trajectories. In column (c) we visualize the predicted pixel displacement between the reference and target timestamp in the Sintel colorization scheme [2].

## REFERENCES

- [1] Z. Teed and J. Deng, “Raft: Recurrent all-pairs field transforms for optical flow,” in *Eur. Conf. Comput. Vis. (ECCV)*, 2020.
- [2] D. J. Butler, J. Wulff, G. B. Stanley, and M. J. Black, “A naturalistic open source movie for optical flow evaluation,” in *Eur. Conf. Comput. Vis. (ECCV)*, 2012.

We compute continuous-time trajectories over similarity

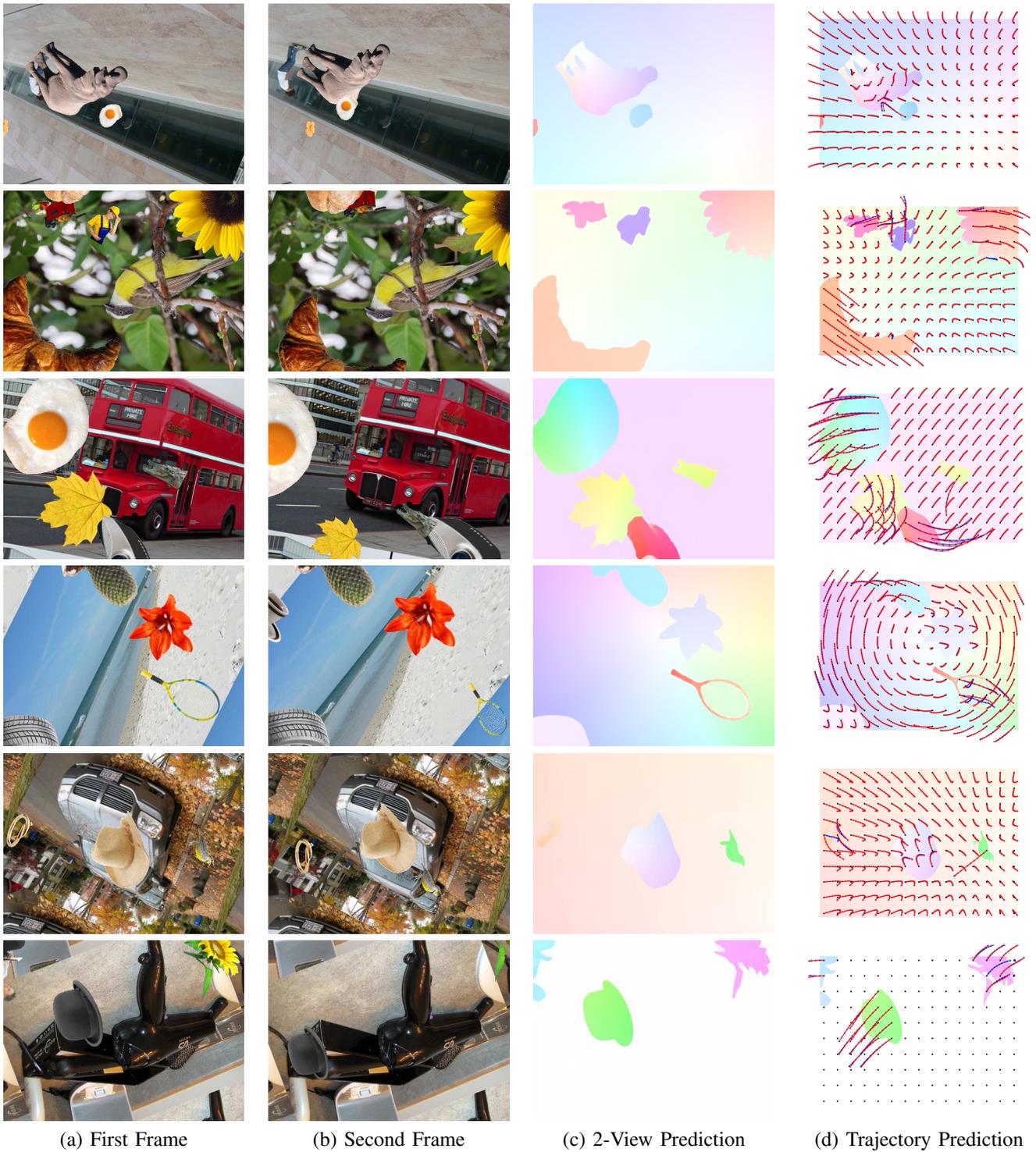


Fig. 3. Predictions of our model (E+I) on our MultiFlow test set. (a) shows the frame at the reference time, (b) the frame at the final timestamp of the trajectory, (c) shows the predicted pixel displacement between reference and final timestamp, and (d) shows the full prediction of the pixel trajectories. The predicted trajectory is shown in blue and the ground truth in red. The background of (d) is the colorization of the 2-view ground truth.