
Appendix

E2PNet: Event to Point Cloud Registration with Spatio-Temporal Representation Learning

Anonymous Author(s)

Affiliation

Address

email

1 A Appendix

2 In this appendix, we first provide the implementation details of our proposed E2PNet and corre-
3 sponding datasets preprocessing (see Sec. A.1). Subsequently, the details of evaluation metrics
4 used for evaluating the generalization of the EP2T module are presented (see Sec. A.2). Finally,
5 we analyze the changes in registration precision under different network settings (see Sec. A.3) and
6 report additional quantitative on other vision-based methods (Event-to-Image Reconstruction, Flow
7 Estimation).

8 A.1 Implementation Details

9 In this section, we first present the implementation details of our proposed Event-Points-to-Tensor
10 (EP2T) module (see Sec. A.1.1). Based on EP2T, we introduce the training strategy of E2PNet and
11 the detailed construction method of event-to-point cloud registration (E2P) dataset (see Sec. A.1.2
12 and Sec. A.1.3). Finally, we introduce the dataset used for generalization experiments (see Sec.
13 A.1.4).

14 A.1.1 Architecture of EP2T

15 **LA and STA.** In the Local Aggregation module (LA), after selecting the neighborhood of each
16 aggregation center, we calculate the spatio-temporal distance between each center and the corre-
17 sponding neighborhood points. For each neighborhood, we use a block consisting of three 2D
18 convolutional with kernel size = $1 * 1$ and batch normalization layers to extract local aggregation
19 distance information. The output dimensions of each layer in the block are (16, 32, 64), and a ReLU
20 activation function is added after the last normalization layer.

21 Subsequently, the features of the aggregation centers will be sent to the Spatio-temporal Separated
22 Attention module (STA). In STA, we use two different attention mechanisms to jointly extract global
23 features [1]. Specifically, we use three different 1D convolution with output channel = 64 and kernel
24 size = 1 to obtain the query matrix, key matrix, and value matrix. To prevent paying too much
25 attention to a certain dimension, we use the multi-head attention ($Head = 4$) mechanism to learn
26 different attention patterns. The difference between self-attention and cross-attention is whether to
27 use another domain to calculate the key matrix.

28 **Tensorized Representation.** After the LA, STA and FP (Feature Propagation) modules, each event
29 point is embedded with high-dimensional (64 channels) spatio-temporal features $f_{FP}(\hat{e}_i)$. To convert
30 point-based features into grid-shaped features, we pioneered the combination of 3D point cloud
31 learning-based and hand-crafted-based methods. Specifically, we separately superimpose the features
32 of each channel to obtain a gridded feature tensor consisting of 64 channels. After that, we perform

33 channel max pooling on all feature points to obtain the global features $f_G(\hat{\mathbf{e}}_i)$ of each point. Finally
 34 we first use three different yet complementary event tensorization methods [2, 3, 4] to convert the
 35 set of $f_G(\hat{\mathbf{e}}_i)$ into different types of 2D grid-shaped sparse feature tensors, and concatenate them to
 36 produce the final tensorized feature map. Specific tensorization methods are as follows:

37 **Sum of Features** [2] separately superimpose the features of each channel to obtain a gridded feature
 38 tensor consisting of 64 channel. The feature value on each polarity and channel can be obtained by
 39 formula (1)

$$\mathbf{F}_{SF(h,w,c)} = \sum_{i=1}^N p_i * k_b(h - h_i) * k_b(w - w_i) * f_{FP}^c(\hat{\mathbf{e}}_i) \quad (1)$$

40 where polarity $p \in \{-1, 1\}$, channel $c \in [1, 64]$, events $\mathbf{e}_i = (h_i, w_i, t_i, p_i)$ and $k_b(x) = \max(0, 1 - |x|)$
 41 is the bilinear sampling kernel, here we use it as an indicator function. Finally, we normalize
 42 them to get a set of grid-shaped feature tensors with shape $64 * H * W$.

43 **Event counting** [2] separates the positive and negative (p_i represents the polarity) events and handles
 44 them separately. The feature value on each polarity channel can be obtained by formula (2)

$$\mathbf{F}_{EC(h,w,p)} = \sum_{i=1}^N k_b(h - h_i) * k_b(w - w_i) * k_b(p - p_i) * f_G(\hat{\mathbf{e}}_i) \quad (2)$$

45 where $p \in \{-1, 1\}$, events $\mathbf{e}_i = (h_i, w_i, t_i, p_i)$ and $k_b(x) = \max(0, 1 - |x|)$ is the bilinear sampling
 46 kernel, here we use it as an indicator function. \mathbf{F}_{EC} superimposes all the features on the same pixel
 47 position. Finally, we normalize them to get a set of grid-shaped feature tensors with shape $2 * H * W$.

48 **Event stacking** [3] divides event clouds into B blocks ($B=3$ in EP2T) equally in time dimension.
 49 Given a set of N input events $\{e_i\}_{i \in [1, N]}$, the events are arranged in the order of appearance time t_i ,
 50 i.e. in ascending order of the value of t_i . N events are re-divided into B blocks according to time t_i ,
 51 the time range of each event stream is $\left[\frac{(b-1)*\Delta t}{B}, \frac{b*\Delta t}{B}\right]$, where $b \in [1, B]$ and $\Delta t = t_N - t_1$. The
 52 calculation method of each blocks is as follow:

$$\mathbf{F}_{ES(h,w,b)} = \sum_{i=1}^N p_i * k_b(h - h_i) * k_b(w - w_i) * f_G(\hat{\mathbf{e}}_i) \quad (3)$$

53 where $k_b(x) = \max(0, 1 - |x|)$. Finally, we normalize \mathbf{F}_{ES} to get a set of grid-shaped feature
 54 tensors with shape $3 * H * W$. Our modification to the method described in the paper [3] involves
 55 multiplying the event eigenvalues by their corresponding positive and negative polarities, and then
 56 accumulating them, rather than determining the overall eigenvalue sign solely based on the polarity of
 57 the latest event. Positive and negative events can cancel each other out, which aligns with the concept
 58 of event polarity as defined by brightness change direction.

59 **Event spatio-temporal Voxelization** [4] not only divides the event point clouds into B blocks ($B=3$ in
 60 EP2T) according to time but also takes the time distance between each event points and the sampling
 61 points as one of the weights of feature aggregation.

$$\mathbf{F}_{EV(h,w,b,p)} = \sum_{i=1}^N k_b(h - h_i) * k_b(w - w_i) * k_b(b - b_i^*) * k_b(p - p_i) * f_G(\hat{\mathbf{e}}_i) \quad (4)$$

62 where $b \in [1, 3]$, $p \in \{-1, 1\}$, $b_i^* = \frac{B*(t_i - t_1)}{t_n - t_i}$ and $k_b(x) = \max(0, 1 - |x|)$. Finally, we normalize
 63 \mathbf{F}_{EV} to get a set of grid-shaped feature tensors with shape $2 * 3 * H * W$. This method not only
 64 encodes the spatial distribution information of events, but also contains the order in which the events
 65 are triggered, giving more granular weights to the spatio-temporal features embedding clouds. In
 66 general, combining the above three event tensorization representation methods, we jointly construct a
 67 sparse tensor with multiple horizons through the coordinate and embedding feature information of
 68 the event clouds.

69 **A.1.2 Training strategy of E2PNet**

70 Our proposed E2PNet is trained with two strategies (tensor-based and point-based). The tensor-based
 71 training strategy is straightforward. After we use EP2T to encode event data into a grid-shaped
 72 tensor, we directly modify the number of input channels (75 in event-based method, 3 for the original
 73 RGB-based method) in the baseline to fit the feature channels of EP2T. Specifically, in DeepI2P
 74 [5], we replace the number of input channels of the first layer in the backbone (resnet34) in the
 75 image branch. In LCD [6], we replace the first convolutional layer of the image branch in the
 76 encoder part. As for the point-based strategy, we add a channel max-pooling layer across the feature
 77 dimension to obtain global features equivalent to the number of channels. To ensure the number of
 78 input channels in subsequent networks is the same as tensor-based methods, we add a linear layer
 79 with output channel = 64 after channel max-pooling. Finally, in both image branch, we replace the
 80 whole encoder backbone of DeepI2P and the encoder part of LCD with this global feature. Other
 81 settings remain unchanged in the original paper [5, 6].

82 **A.1.3 Datasets Preprocessing for E2PNet**

83 Since there is no dataset for the E2P task, we propose to use multi-sensor SLAM datasets to construct
 84 the registration relationship GT between event cameras and 3D point clouds. We select MVSEC [7]
 85 and VECtor [8] as these datasets use LiDAR, traditional cameras and event cameras simultaneously,
 86 and have good calibration. With these calibration parameters and pose GT, we can establish the GT
 87 matching relationship between events (images) and point clouds for E2P tasks through the camera
 88 projection model. (see Eq. 5).

$$\mathbf{Z} \begin{pmatrix} u \\ v \\ 1 \end{pmatrix} = \begin{pmatrix} f_x & 0 & c_x \\ 0 & f_y & c_y \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} X \\ Y \\ Z \end{pmatrix} = \mathbf{K} \mathbf{P}_c = \mathbf{K} \mathbf{T}_{c_w} \mathbf{P}_w \quad (5)$$

89 where K is the internal parameter of the camera. We first use LiDAR mapping algorithm [9] to
 90 construct the complete point cloud map \mathbf{P}_w , and then use GT pose \mathbf{T}_{c_w} to transform the point cloud
 91 map into the camera coordinate system to obtain \mathbf{P}_c .

92 Since the event camera has the same optical lens principle as the traditional camera, we can obtain
 93 the point cloud within the viewing frustum of the event camera through the camera projection model
 94 (Eq. 5). However, camera projection for all point clouds is computationally expensive. Equivalently,
 95 we construct a camera model in the world coordinate system and use this camera model to construct
 96 a quadrangular viewing frustum. By applying four sets of plane equations, we can filter the range
 97 of point cloud coordinates effectively. Each screening operation significantly reduces the number
 98 of point clouds that require future calculations. Experimental results show that the calculation time
 99 is reduced by more than half when selecting points in the quadrilateral viewing frustum compared
 100 with projecting all points and performing visibility filtering. The detailed experimental setup in two
 101 different datasets is as follows:

102 **MVSEC** [7] is a well-known event-based dataset. MVSEC provide data on various motion mode
 103 (carried on a handheld rig, flown by a hexacopter, driven on top of a car and mounted on a motorcycle)
 104 in various scenarios. Since only indoor scenes provide high-precision pose and high-quality point
 105 cloud data, we only use indoor scene data for E2P tasks. We use the indoor- x and indoor- y sequences
 106 for training and testing respectively, where $x \in [1, 3]$ and $y = 4$. This experimental setup would
 107 generate a total number of 20,400 event (image) and point cloud pairs with the corresponding pose
 108 for training and 1,610 pairs for testing. In addition, each point cloud is augmented by adding random
 109 rotation (up to 0.5 degrees) and translation (up to 0.01 m).

110 **VEctor** [8] is the first event-based SLAM benchmark dataset captured by a full hardware-
 111 synchronized sensor suite. We use the VECtor dataset mainly to test the E2P effect in large-scale
 112 indoor scenarios. The VECtor dataset contains three different scenes of campus building interiors,
 113 each captured by two different motion modalities. We use the units-dolly, units-scooter, corridors-
 114 dolly and corridors-walk sequences for training, and the school-dolly and school-scooter sequences
 115 for evaluation. This experimental setup would generate 3,025 pairs for training and 1,544 pairs for
 116 testing. It is worth noting that our experimental setup uses completely different scenarios when
 117 training and testing, which will verify the generalization of our EP2T method to new scenarios. The
 118 same data augmentation method as MVSEC is also used.

119 A.1.4 Datasets Preprocessing for Generalization Experiments

120 **MVSEC** [7] is also evaluated for optical flow estimation. However, the optical flow GT is sparse
121 due to using the pose and depth map provided by LiDAR. Follow [2, 3], we use the *outdoor-day-2*
122 sequence (more than 12K grayscale frames) for training and test it on the *outdoor-day-1* sequences.
123 During training, we randomly select a processing window that spans 1/45 second, from which we
124 extract all events and the optical flow GT. However, the training data will be discarded if the number
125 of events in the window less than 10,000. It should be noted that the starting position of the processing
126 window is chosen randomly, so it may not be aligned with the timestamp of the optical flow GT.
127 To address this special condition, an interpolation operation may be necessary. For event-to-image
128 reconstruction, we use the same sequence as flow estimation for training and testing.

129 **ECD** [10] contains a set of the asynchronous event stream, intensity images at about 24Hz, GT
130 camera poses from a motion-capture system with sub-millimeter precision at 200Hz. As introduced
131 in BTEB [11], we use the same sequences cut for testing. Unlike the MVSEC dataset, since BTEB
132 does not require GT reconstructed images as a supervisory signal, we randomly but continuously
133 sample the entire event sequence ($N = 8,192$) during training. During testing, we extract all events
134 between image frames, and after random sampling, we use these sparse event data to reconstruct an
135 image.

136 **N-Caltech101 and N-CARS** [12, 13] are two large public datasets for event-based classification.
137 N-Caltech101 is an event camera version of the Caltech101 [12] dataset, which was created by
138 moving an event camera that focused on an LCD monitor displaying the original Caltech101 data.
139 Following the example of the original paper [12], we randomly select 15 samples from each class
140 for testing. In N-CARS dataset, an event camera was mounted behind the windshield of a moving
141 car during collection, and each sample contains 59,249 events. The whole dataset comprises 12,336
142 car samples and 11,693 background samples, which divided 7,939 cars and 7,482 backgrounds for
143 training and others for testing. Similarly, we sample 8,192 events from a large number of events in
144 each sample.

145 A.2 Evaluation Metrics

146 In order to evaluate the generalization of our proposed EP2T module, we employ multiple evaluation
147 metrics [2, 14, 15] to compare the results of different tasks.

148 **Flow Estimation.** Average end-point error (AEE [2]) denotes the distance between the end-points
149 of the predicted (Y') and the GT (Y) flow vectors:

$$AEE = \sum_{x,y} \left\| \begin{pmatrix} u(x,y)_{Y'} \\ v(x,y)_{Y'} \end{pmatrix} - \begin{pmatrix} u(x,y)_Y \\ v(x,y)_Y \end{pmatrix} \right\|_2 \quad (6)$$

150 where u, v represent the horizontal and vertical optical flow value respectively. Follow [2, 3], we
151 limit the computation of AEE to pixels in which at least one event was observed. In addition to
152 pixel-level evaluation, we also perform a global outlier ratio analysis. Our experiments regard pixels
153 with $AEE > 3$ as outliers.

$$Outlier = \frac{\sum_{i=1}^{H \times W} \mathbb{1}(AEE_i > 3)}{H \times W} \times 100\% \quad (7)$$

154 where H and W denotes the solution of the flow image, $\mathbb{1}(\cdot)$ is the indicator function, which equals 1
155 if the AEE of pixel i is greater than 3, and 0 otherwise.

156 **Event-to-Image Reconstruction.** Mean squared error (MSE) is a classic evaluation metric to
157 compare the distance of two different images. Unlike AEE (defined in Eq. 6), the reconstruction task
158 is to evaluate all pixels of the image (regardless of whether event exists).

$$MSE = \frac{1}{n} \sum_{i=1}^{H \times W} \left\| \vec{Y}_i - \vec{Y}'_i \right\|_2 \quad (8)$$

159 where Y is the GT grayscale image and Y' is the reconstructed image by event-based methods.
160 Similarly, in addition to pixel-level evaluation, reconstruction tasks often focus on the semantic

Table 1: **Registration performance (DeepI2P as baseline in MVSEC-E2P dataset) under different sampling methods.**

Method	Uniform	Random	Farthest Point Sampling [17]	Voxel	Surface Event Sampling [18]
RE(°)(↓)	5.127	12.215	6.521	8.622	5.095
TE(m)(↓)	0.164	4.200	0.249	0.212	0.166
Time(ms)(↓)	1.1	0.3	796	8	421

161 similarity (LPIPS [15]) and structural similarity (SSIM [14]) between the reconstructed image and
 162 the GT image to prevent perceptual differences. LPIPS is an image quality evaluation index based on
 163 deep learning, which measures the semantic similarity between two images, defined as follows:

$$\text{LPIPS} = \frac{1}{N} \sum_{i=1}^N D_{\text{net}}(Y'_i, Y_i) \quad (9)$$

164 where D_{net} represents the distance or difference metric between image patches computed by a
 165 pre-trained deep network (AlexNet [16] in experiments), N is the number of patch while the patch
 166 size is 32×32 . SSIM measures how similar two images are in terms of structure, brightness, and
 167 contrast, detailed defined below:

$$\text{SSIM} = \frac{(2\mu_{Y'}\mu_Y + C_1)(2\sigma_{Y'Y} + C_2)}{(\mu_{Y'}^2 + \mu_Y^2 + C_1)(\sigma_{Y'}^2 + \sigma_Y^2 + C_2)} \quad (10)$$

$$C_1, C_2 = (k_1L)^2, (k_2L)^2$$

168 where μ denotes the mean of pixel values, σ denotes the standard deviation of pixel values, and $\sigma_{Y'Y}$
 169 indicates the covariance between pixel values of two images. C_1 and C_2 are constants for stable
 170 calculations, L is the dynamic range of pixel values (255 in common) and $k_1 = 0.01, k_2 = 0.03$.

171 **Object Recognition.** Accuracy is a commonly used evaluation metric for classification model per-
 172 formance, which is used to measure the proportion of samples correctly classified by the classification
 173 model in the prediction process. The simple definition is as follows:

$$\text{Accuracy} = \frac{M}{N} \times 100\% \quad (11)$$

174 where M indicates the number of samples correctly classified by the classification model in the
 175 prediction, N represents the total number of predicted samples.

176 A.3 Additional Results

177 **The impacts of sampling methods on performance.** To better demonstrate the performance
 178 of E2P under different aggregation centers, we further report the results under different sampling
 179 methods. As shown in Tab. 1, we compare the performance between Uniform Sampling (US),
 180 Random Sampling (RS), Farthest Point Sampling (FPS [17]), Voxel Sampling (VS) and Surface
 181 Event Sampling (SES [18]). For a fair comparison, we choose DeepI2P as the baseline and evaluate
 182 in MVSEC-E2P dataset. Experiments show that the simple uniform sampling method used in EP2T
 183 is relatively excellent regarding time efficiency, only slightly slower than the spatio-temporal random
 184 sampling method. In addition, because uniform sampling overcomes the inhomogeneity of local
 185 redundancy and global sparseness of event data, the registration precision also shows good potential,
 186 which is comparable to the learning-based event surface sampling method.

187 **Qualitative examples.** Here, we provide visualizations in Fig. 1 and Fig. 2 to present the effect
 188 of our EP2T module under different vision-based tasks (optical flow estimation, event-to-image
 189 reconstruction). Specifically, we demonstrate the event-to-image reconstruction effect of the EP2T
 190 module on the ECD [10] dataset with BTEB [11] as the baseline. As the limitation we described, it
 191 would be unwise to take all events between image frames for reconstruction like most frameworks but
 192 randomly sample 8,192 events. From qualitative experiments, it can be seen that even if the sampling
 193 operation reduces global information, our framework still maintains good reconstruction performance
 194 (see Fig. 1). Similarly, we also demonstrated the optical flow estimation effect of the cGAN [19]
 195 network as a baseline under the MVSEC [7] dataset. Dense optical flow denotes flow with all pixels,
 196 while sparse optical flow denotes masked flow at the pixels with events (see Fig. 2).

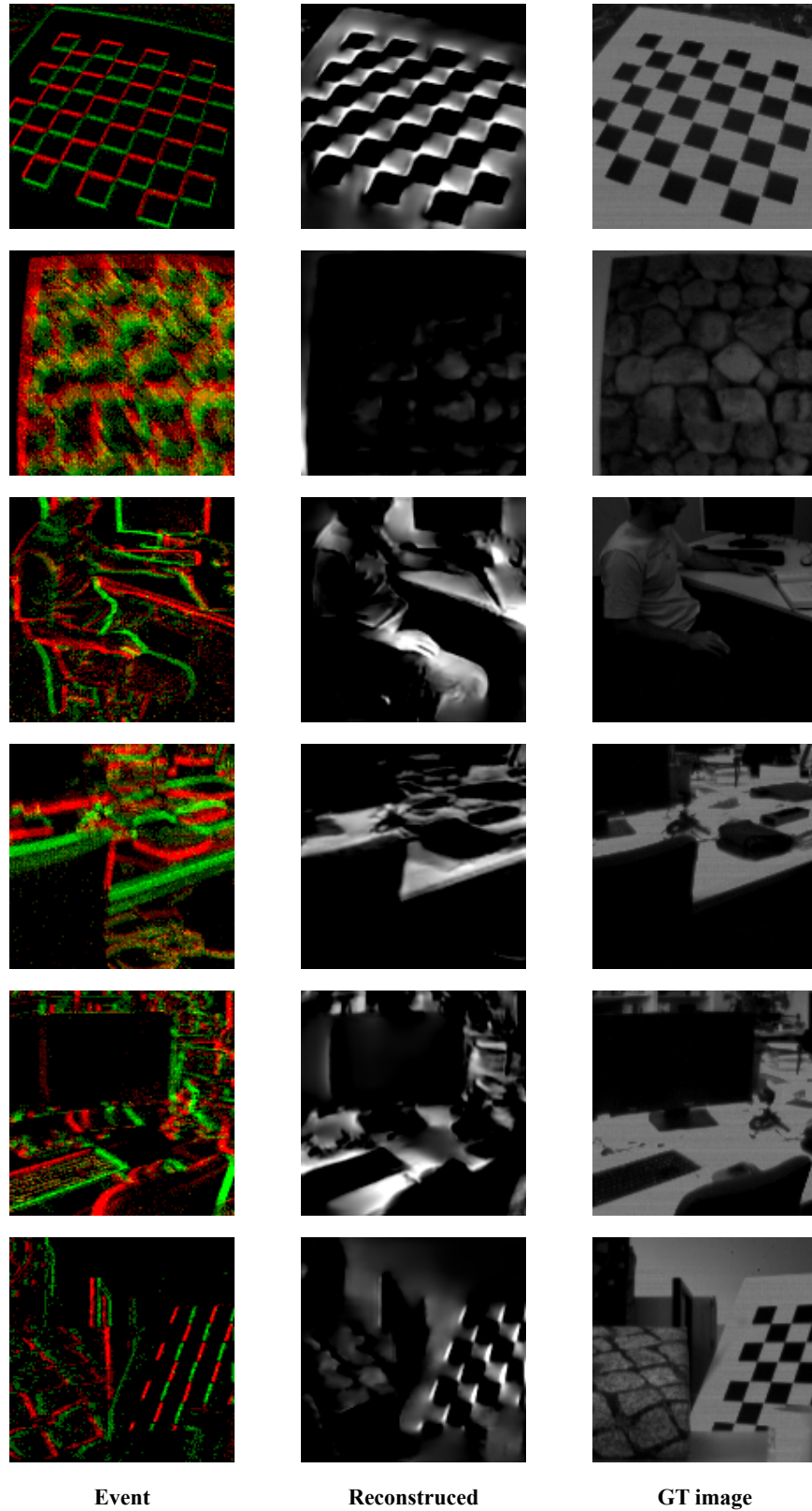


Figure 1: Example of event-to-image reconstruction.

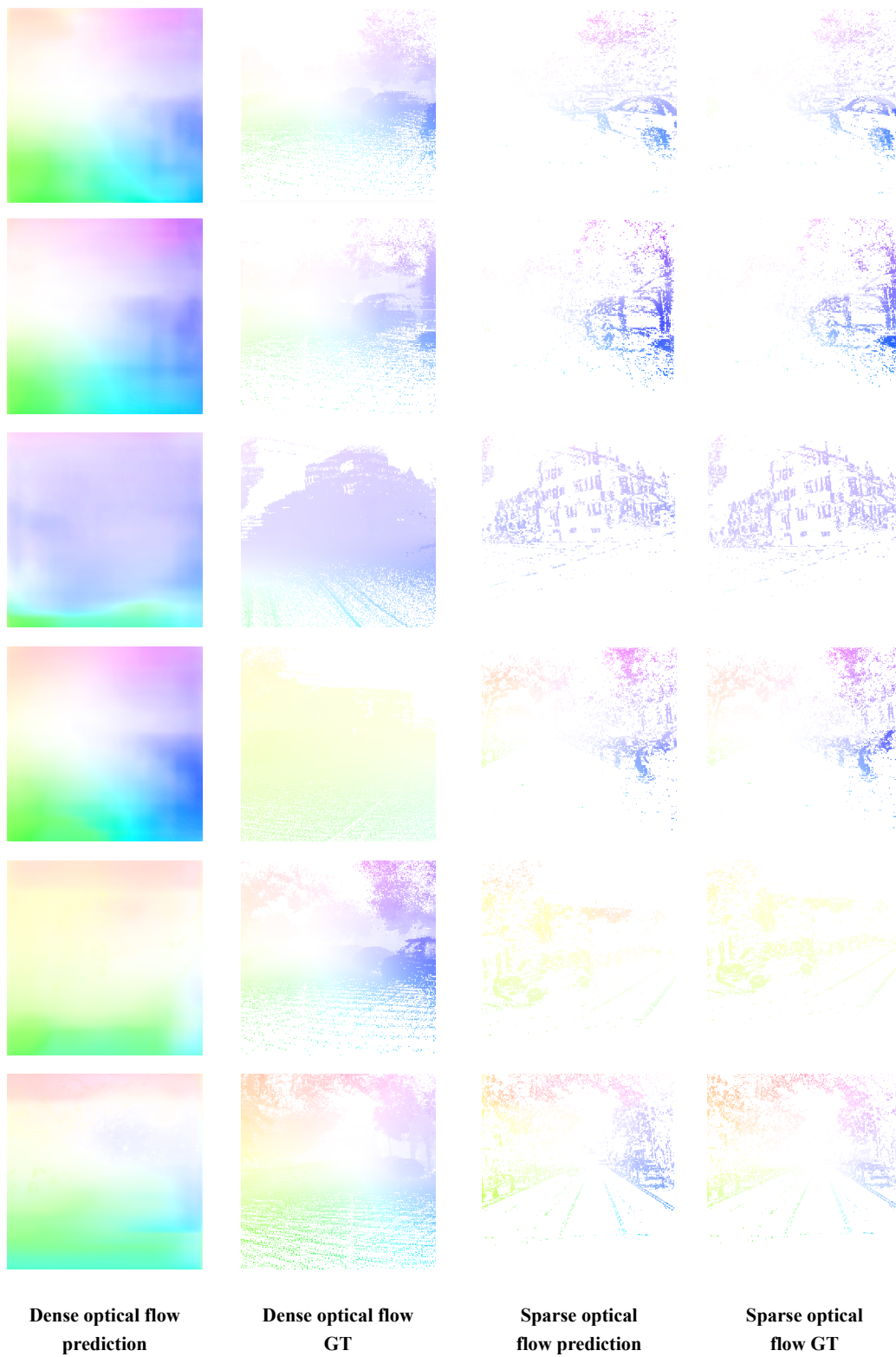


Figure 2: **Example of optical flow estimation.**

References

- 197
- 198 [1] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
199 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information*
200 *processing systems*, 30, 2017.
- 201 [2] Alex Zihao Zhu, Liangzhe Yuan, Kenneth Chaney, and Kostas Daniilidis. Ev-flownet: Self-
202 supervised optical flow estimation for event-based cameras. *arXiv preprint arXiv:1802.06898*,
203 2018.
- 204 [3] Mohammad Mostafavi, Lin Wang, and Kuk-Jin Yoon. Learning to reconstruct hdr images
205 from events, with applications to depth and flow prediction. *International Journal of Computer*
206 *Vision*, 129:900–920, 2021.
- 207 [4] Alex Zihao Zhu, Liangzhe Yuan, Kenneth Chaney, and Kostas Daniilidis. Unsupervised
208 event-based learning of optical flow, depth, and egomotion. In *Proceedings of the IEEE/CVF*
209 *Conference on Computer Vision and Pattern Recognition*, pages 989–997, 2019.
- 210 [5] Jiaxin Li and Gim Hee Lee. Deepi2p: Image-to-point cloud registration via deep classification.
211 In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*,
212 pages 15960–15969, 2021.
- 213 [6] Quang-Hieu Pham, Mikaela Angelina Uy, Binh-Son Hua, Duc Thanh Nguyen, Gemma Roig,
214 and Sai-Kit Yeung. Lcd: Learned cross-domain descriptors for 2d-3d matching. In *Proceedings*
215 *of the AAAI Conference on Artificial Intelligence*, volume 34, pages 11856–11864, 2020.
- 216 [7] Alex Zihao Zhu, Dinesh Thakur, Tolga Özaslan, Bernd Pfrommer, Vijay Kumar, and Kostas
217 Daniilidis. The multivehicle stereo event camera dataset: An event camera dataset for 3d
218 perception. *IEEE Robotics and Automation Letters*, 3(3):2032–2039, 2018.
- 219 [8] Ling Gao, Yuxuan Liang, Jiaqi Yang, Shaoxun Wu, Chenyu Wang, Jiaben Chen, and Laurent
220 Kneip. Vector: A versatile event-centric benchmark for multi-sensor slam. *IEEE Robotics and*
221 *Automation Letters*, 7(3):8217–8224, 2022.
- 222 [9] David Levin. The approximation power of moving least-squares. *Mathematics of computation*,
223 67(224):1517–1531, 1998.
- 224 [10] Elias Mueggler, Henri Rebecq, Guillermo Gallego, Tobi Delbruck, and Davide Scaramuzza.
225 The event-camera dataset and simulator: Event-based data for pose estimation, visual odometry,
226 and slam. *The International Journal of Robotics Research*, 36(2):142–149, 2017.
- 227 [11] Federico Paredes-Vallés and Guido CHE de Croon. Back to event basics: Self-supervised
228 learning of image reconstruction for event cameras via photometric constancy. In *Proceedings*
229 *of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3446–3455,
230 2021.
- 231 [12] Garrick Orchard, Ajinkya Jayawant, Gregory K Cohen, and Nitish Thakor. Converting static
232 image datasets to spiking neuromorphic datasets using saccades. *Frontiers in neuroscience*,
233 9:437, 2015.
- 234 [13] Amos Sironi, Manuele Brambilla, Nicolas Bourdis, Xavier Lagorce, and Ryad Benosman.
235 Hats: Histograms of averaged time surfaces for robust event-based object classification. In
236 *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1731–
237 1740, 2018.
- 238 [14] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment:
239 from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4):600–
240 612, 2004.
- 241 [15] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unrea-
242 sonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE*
243 *conference on computer vision and pattern recognition*, pages 586–595, 2018.

- 244 [16] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep
245 convolutional neural networks. *Communications of the ACM*, 60(6):84–90, 2017.
- 246 [17] Carsten Moenning and Neil A Dodgson. Intrinsic point cloud simplification. *Proc. 14th*
247 *GrahiCon*, 14:23, 2004.
- 248 [18] Zhiwen Chen, Jinjian Wu, Junhui Hou, Leida Li, Weisheng Dong, and Guangming Shi. Ecsnet:
249 Spatio-temporal feature learning for event camera. *IEEE Transactions on Circuits and Systems*
250 *for Video Technology*, 2022.
- 251 [19] Mohammad Mostafavi, Lin Wang, and Kuk-Jin Yoon. Learning to reconstruct hdr images
252 from events, with applications to depth and flow prediction. *International Journal of Computer*
253 *Vision*, 129:900–920, 2021.