# AW-GATCN: Adaptive Weighted Graph Attention Convolutional Network for Event Camera Data Joint Denoising and Object Recognition

# Haiyu Li

School of Electrical and Electronic Engineering
The University of Sheffield
Sheffield, United Kingdom
hli108@sheffield.ac.uk

# Charith Abhayaratne

School of Electrical and Electronic Engineering
Centre for Machine Intelligence
The University of Sheffield
Sheffield, United Kingdom
c.abhayaratne@sheffield.ac.uk

Abstract—Event cameras, which capture brightness changes with high temporal resolution, inherently generate a significant amount of redundant and noisy data beyond essential object structures. The primary challenge in event-based object recognition lies in effectively removing this noise without losing critical spatial-temporal information. To address this, we propose an Adaptive Graph-based Noisy Data Removal framework for Event-based Object Recognition. Specifically, our approach integrates adaptive event segmentation based on normalized density analysis, a multifactorial edge-weighting mechanism, and adaptive graph-based denoising strategies. These innovations significantly enhance the integration of spatiotemporal information, effectively filtering noise while preserving critical structural features for robust recognition. Experimental evaluations on four challenging datasets demonstrate that our method achieves superior recognition accuracies of 83.77%, 76.79%, 99.30%, and 96.89%, surpassing existing graph-based methods by up to 8.79%, and improving noise reduction performance by up to 19.57%, with an additional accuracy gain of 6.26% compared to traditional Euclidean-based techniques.

Index Terms—event camera, denoising, GATCN, object recognition

### I. INTRODUCTION

The advent of artificial intelligence has propelled computer vision to the forefront of real-world applications such as autonomous driving, drone navigation, and surveillance, where swift and accurate object recognition is paramount. Traditional video cameras, limited by low frame rates and excessive data redundancy, fall short in these dynamic settings. Although high-speed cameras capture over 1,000 frames per second, their high cost limits practicality. In contrast, event cameras [1], which only record changes in scene brightness, significantly minimize data redundancy and are unaffected by motion blur. They offer microsecond-level temporal resolution and low latency, making them ideal for environments that demand rapid and reliable data processing.

Unlike conventional cameras that output continuous twodimensional images, event cameras are triggered by significant changes in pixel brightness, efficiently eliminating most of

<sup>0</sup>This is the preprint version of a paper accepted at IJCNN 2025. The final version will appear in the IJCNN 2025 proceedings.

irrelevant background information. However, the asynchronous and sparse data they generate pose significant challenges for traditional frame-based processing techniques [2]. Researchers typically convert event streams into 2D frames or 3D voxel grids [3]–[5], a process that compromises the data's inherent sparsity and temporal resolution, leading to potential information loss. The absence of a standardized conversion method further complicates data processing, as application-specific needs require customized approaches, yielding inconsistent results across different scenarios.

To fully leverage the unique characteristics of event cameras—namely, the sparsity and asynchrony of event data—researchers have explored innovative processing methodologies such as temporal surface-based methods [6], [7] and spiking neural networks (SNNs) [8], [9]. These approaches, which process data on an event-by-event basis, are designed to maintain low latency. However, their efficacy in complex tasks can be limited due to the sensitivity to parameter settings and the intricacies of their training processes. In response to these challenges, recent advancements have introduced a more efficient strategy involving compact graph representations [10]-[12]. These methods model event sequences as graphs within a cloud of events using graph convolutional networks (GCNs), which have achieved state-of-the-art performance. Despite their successes, these graph-based approaches primarily rely on a simplistic radius-based noise management strategy-connecting nodes only if they are within a predetermined Euclidean distance. This technique often proves inadequate for effectively handling noise and lacks the flexibility needed for adapting to dynamically changing environments.

Based on these observations, we propose an adaptive graph formulation-based noise reduction algorithm integrated with a graph convolutional neural network (GCN) that incorporates attention mechanisms, enabling efficient and accurate processing of event data. Traditional radius-based methods rely on fixed Euclidean distances, limiting adaptability and robustness while overlooking other informative graph features.

Our approach overcomes these limitations by incorporating multilevel weights, dynamically adjusting weight thresholds and leveraging a graph attention mechanism for enhanced feature aggregation and classification.

The main contributions of this paper are as follows:

- Multi-Factor Edge Weighting: A robust edge weighting mechanism that incorporates Euclidean distance, velocity, angular difference, and polarity consistency, ensuring accurate modeling of event point relationships, even in noisy environments.
- Adaptive graph formulation-based Noise Reduction:
   A dynamic noise reduction strategy that adapts the formulation of the underlying graph of event based on the variance in node distribution, effectively filtering out noisy events by preserving the key event data in sparse areas and removing excess in dense regions.
- Graph Attention Convolutional Network: A GCN guided by multi-factor edge weighting, selectively emphasizing relevant neighboring features and achieving enhanced data representation and object recognition accuracy.

#### II. RELATED WORK

Current methods for event data processing can be broadly categorized into frame-based conversion methods, graph-based methods, and deep learning methods. Frame-based methods, such as those proposed by [3], [4], convert event data into pseudo-images, but this often leads to a loss of temporal resolution. Graph-based methods [13], [14] maintain the sparsity of event data but face challenges in processing efficiency. Deep learning methods, particularly those utilizing Graph Attention Neural Networks [15], are adept at handling complex graph-structured data but often struggle with noisy data [16].

One approach for event cameras is to use Spiking Neural Networks (SNN) [17], again a biologically inspired design. SNNs exploit the sparsity and asynchronous nature of event data, but due to their non-microscopic nature, training such networks is very difficult. To improve the temporal resolution, Zhu et al [18], [19] suggested discretizing the temporal dimension into consecutive time segments and accumulating the events into a voxel grid by linear weighted accumulation similar to bilinear interpolation. Messikommer et al [20] further exploited spatial and temporal sparsity by employing sparse convolution [21] and developing recursive convolution formulations. However, they still operate on sparse volumes and 3D convolution is computationally expensive for dealing with large event clouds.

Recent studies, e.g., further use a framework similar to PointNet [22], [23], which utilizes a multilayer perceptron (MLP) to learn the features of each point separately and then outputs object-level responses (e.g., categorical labels) via a global max operation. For event processing, Sekikawa et al [24] developed for the first time a recursive architecture called EventNet. Specifically, it recursively represents the dependencies of causal events on outputs through a new temporal encoding and aggregation scheme and pre-computes

# Algorithm 1 Adaptive Event Point Segmentation

1:  $x_{\text{range}} \leftarrow x_{\text{max}} - x_{\text{min}}$ 

```
2: y_{\text{range}} \leftarrow y_{\text{max}} - y_{\text{min}}
3: t_{\text{range}} \leftarrow t_{\text{max}} - t_{\text{min}}
4: \rho_{\text{norm}} \leftarrow \frac{N_{\text{points}}}{x_{\text{range}} \times y_{\text{range}} \times t_{\text{range}}}
5: N_{\text{window}} \leftarrow \max(N_{\text{min}}, \rho_{\text{norm}} \times C_{\text{scale}})
 6: N_{\text{windows}} \leftarrow \left\lceil \frac{N_{\text{points}}}{N_{\text{window}}} \right\rceil
 7: for w \leftarrow 1 to N_{\text{windows}} do
                 Sort data in window w by t
 8:
 9:
                 x_{\text{span}} \leftarrow x_{\text{max}} - x_{\text{min}}
10:
                 y_{\text{span}} \leftarrow y_{\text{max}} - y_{\text{min}}
                 t_{\text{span}} \leftarrow t_{\text{max}} - t_{\text{min}}
11:
                 N_{\text{voxels}} \leftarrow \sqrt{x_{\text{span}} \times y_{\text{span}} \times t_{\text{span}}}
12:
13:
                 N_{\text{voxels}} \leftarrow \max(N_{\text{min\_vox}}, \min(N_{\text{max\_vox}}, N_{\text{voxels}}))
                 Divide window w into N_{\text{voxels}} based on time intervals
14:
                 for m \leftarrow 1 to N_{\text{voxels}} do
15:
                          Extract event points in voxel m
16:
                 end for
17:
18: end for
```

the features of nodes that correspond to particular spatial coordinates and polarities.

# III. METHODOLOGY

Processing event videos from event cameras requires managing large volumes of noisy event points [25], making direct processing computationally demanding. To address this, we propose an adaptive segmentation algorithm that first divides the input data into balanced windows based on normalized density and then subdivides each window into voxels using the square root law to balance temporal and spatial dimensions. For noise reduction, inspired by [26], we employ an adaptive algorithm that dynamically adjusts the weighting radius based on multiple event point features, filtering out noise. These weights are then integrated with a graph attention mechanism to selectively focus on relevant neighboring features, improving object recognition performance. An overview of the framework is shown in Figure 1.

#### A. Adaptive Event Point Segmentation

Event data generated by event cameras is represented as a point cloud, denoted as  $\mathcal{E} = \{(x_k, y_k, t_k, p_k)\}_{k=1}^{N_{\text{points}}}$ , where each point k includes spatial coordinates  $(x_k, y_k)$ , a timestamp  $t_k$ , and a polarity  $p_k$ . Here,  $N_{\text{points}}$  represents the total number of event points. Algorithm 1 outlines the preprocessing and segmentation procedures, which are conducted as follows:

1) Normalized Density: To effectively segment event-based data, we commence by calculating the normalized density of data points. This calculation begins with determining the spatial range, capturing the data's extent in the x and y dimensions. By evaluating the distribution of data points along these axes, we establish the spatial boundaries essential for understanding the overall spatial distribution.

We then assess the temporal range of the data, reflecting the duration over which the events are recorded. This temporal

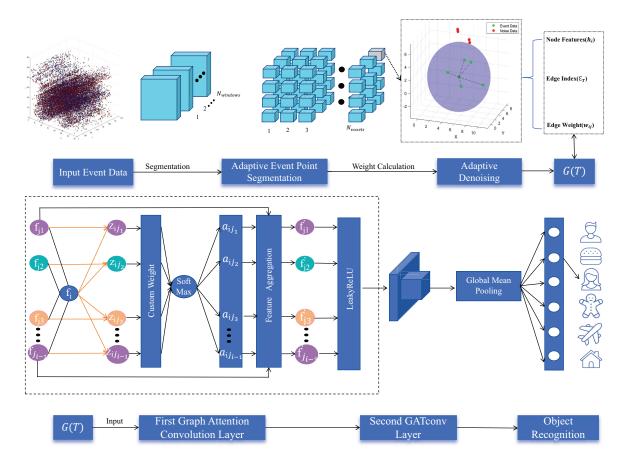


Fig. 1: The adaptive edge weight  $w_{ij}$ , designed to capture event point relevance, facilitates noise filtering and attention-based feature aggregation, enhancing robust recognition by emphasizing the most informative node connections.

assessment aids in comprehending the temporal dynamics of the event stream, enhancing our understanding of its temporal characteristics.

With both spatial and temporal ranges defined, the normalized density is computed as the number of data points per unit volume, encompassing both spatial and temporal dimensions. The formula for normalized density is given by:

$$\rho_{\rm nor} = \frac{N_{\rm points}}{(x_{\rm max} - x_{\rm min}) \cdot (y_{\rm max} - y_{\rm min}) \cdot (t_{\rm max} - t_{\rm min})}, \quad (1)$$

where  $N_{\rm points}$  denotes the total number of data points,  $(x_{\rm max}, x_{\rm min}, y_{\rm max}, y_{\rm min})$  represent the spatial boundaries, and  $t_{\rm max}, t_{\rm min}$  are the temporal boundaries.

2) Determination of Adaptive Window Size: To ensure an even distribution of data points across windows, we adaptively adjust the window size based on the normalized density. The number of points per window is determined by the following heuristic:

$$n_{\text{window}} = \max(N_{\text{min}}, \rho_{\text{normalized}} \cdot C_{\text{scale}}),$$
 (2)

where  $N_{\rm min}$  is the minimum number of points per window, ensuring that windows are not overly sparse, and  $C_{\rm scale}$  is a scaling factor to control the overall window size. This adjustment guarantees that each window contains a sufficient

number of points for robust analysis, mitigating issues such as noise interference or sparse data regions.

After determining the optimal number of points per window, the total number of windows required to cover all data points is calculated as:

$$N_{\text{windows}} = \left\lceil \frac{N_{\text{points}}}{n_{\text{window}}} \right\rceil, \tag{3}$$

where  $\lceil \cdot \rceil$  denotes the ceiling function, ensuring complete coverage of all data points across the windows.

This balanced distribution of data points within each window supports adaptive segmentation and facilitates reliable processing across varying data densities. The event point set  $\mathcal{E}$  is divided into multiple windows  $\{\mathcal{W}_w\}_{w=1}^{N_{\text{windows}}}$ .

3) Determination of Adaptive Voxel Count: To further refine segmentation within each window, the voxel count is determined based on the **square root law**, which balances both spatial and temporal spans. The square root law is a heuristic often used in multidimensional systems to proportionally balance different dimensions by taking the square root of their product. In our approach, this allows us to determine an optimal voxel count that preserves data resolution while maintaining computational efficiency.

The number of voxels within each window is calculated as:

$$N_{\text{voxels}} = \sqrt{(X) \cdot (Y) \cdot (T)},$$
 (4)

where  $X=x_{\rm max}-x_{\rm min},\,Y=y_{\rm max}-y_{\rm min},\,{\rm and}\,\,T=t_{\rm max}-t_{\rm min}$  represent the spatial and temporal ranges of the data within each window. This formula ensures that the voxelization process captures essential data characteristics. By applying the square root law, we achieve a segmentation that is both efficient and sufficiently detailed for subsequent analysis.

After dividing each window  $\mathcal{W}_w$  into multiple voxels  $\{V_m\}_{m=1}^{N_{\text{voxels}}}$ , each voxel  $V_m$  contains a subset of event points characterized by their spatial, temporal, and polarity attributes. This segmentation enables each voxel  $V_m$  to encapsulate a localized subset of events, represented as follows:

$$V_m = \{(x_k, y_k, t_k, p_k)\}. \tag{5}$$

Here, k represents the index of event points within voxel  $V_m$ , and its range is determined by the number of points that fall within  $V_m$ . Specifically, if  $V_m$  contains  $N_{\text{points},m}$  event points, then  $k=1,2,\ldots,N_{\text{points},m}$ .

# B. Adaptive Denoising

AW-GATCN applies a weight-based noise reduction method that uses the variance of the normalized degree matrix to improve 3D structure recognition. This approach considers Euclidean distance, angular velocity difference, velocity magnitude difference, and polarity consistency between nodes. The optimal weight threshold is adaptively determined by maximizing the variance in node distribution across graphs within each voxel, unlike traditional methods with manually set thresholds. To maintain clarity and follow common graph processing conventions, we use k to denote event point indices during segmentation, switching to i and j for node indices in subsequent graph-based steps.

1) Custom Weight Calculation: For each pair of event points i and j within voxel  $V_m$ , the edge weight  $w_{ij}$  is computed as:

$$w_{ij} = \alpha \cdot D_{ij} + \beta \cdot \Delta v_{ij} + \gamma \cdot \theta_{ij} + \delta \cdot P_{ij}, \tag{6}$$

- $D_{ij}$ : The Euclidean distance between nodes i and j, representing the spatial distance.
- $\Delta v_{ij}$ : The magnitude difference between the velocity vectors of nodes i and j.
- $\theta_{ij}$ : The angle difference between the velocity vectors or planar vectors of nodes i and j.
- $P_{ij}$ : Polarity consistency, a binary indicator representing whether the polarities of nodes i and j are consistent (0 if consistent, 1 if inconsistent).

The velocity vector  $\mathbf{v}_i$  for each event point is computed based on the spatial and temporal differences between neighboring event points. For two event points  $e_i = (x_i, y_i, t_i)$  and  $e_j = (x_j, y_j, t_j)$ , the velocity vector is defined as:

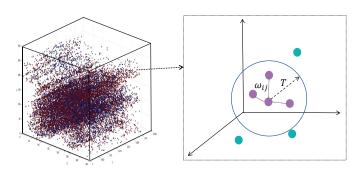


Fig. 2: The circular region represents the optimal threshold T determined by maximizing variance. Purple nodes indicate high-correlation points with weights  $w_{ij}$  less than T, while teal nodes are excluded.

$$\mathbf{v}_i = \left(\frac{x_j - x_i}{t_j - t_i}, \frac{y_j - y_i}{t_j - t_i}\right). \tag{7}$$

This vector represents the "movement" of event point  $e_i$  through space and time. When  $t_i = t_j$ , only the angular and magnitude differences in the 2D plane velocity vectors are calculated; if  $t_i \neq t_j$ , both spatial and temporal velocity differences are computed. Since polarity does not influence velocity vector calculations, it is omitted here.

The angular difference between two velocity vectors describes their variation in motion direction and is given by:

$$\cos \theta_{ij} = \frac{\mathbf{v}_i \cdot \mathbf{v}_j}{|\mathbf{v}_i||\mathbf{v}_j|},\tag{8}$$

where  $\mathbf{v}_i \cdot \mathbf{v}_j$  denotes the dot product, and  $|\mathbf{v}_i|$ ,  $|\mathbf{v}_j|$  are the magnitudes of the vectors. A smaller angle  $\theta_{ij}$  implies similar motion directions, while a larger angle indicates more pronounced directional differences.

In the weight calculation (see Equation (6)), the angular difference between velocity vectors plays a crucial role, as  $\theta_{ij}$  captures the similarity in motion directions. This term aids in identifying local motion patterns, such as sliding or rotating edges. By incorporating angular differences into the weight function, the method more accurately captures dynamic relationships between event points, thereby enhancing overall performance.

After conducting experiments, we finalized the weight parameters  $\alpha=0.7$ ,  $\beta=0.1$ ,  $\gamma=0.1$ , and  $\delta=0.1$ , as these values demonstrated effective performance in capturing dynamic relationships between event points (see Section IV-C1).

2) Dynamic Threshold Adjustment Based on Normalized Degree Distribution: Each voxel contains both event points and noise. To establish full connectivity for noise reduction, we first apply the Minimum Spanning Tree (MST) algorithm to determine the minimum weight that connects all event points within the voxel. This weight serves as the upper limit for the noise reduction threshold, ensuring that all points are fully connected. To effectively filter out noise while preserving

# **Algorithm 2** Computing Optimal Threshold T

```
1: for n = 0 to m - 1 do
                                                 \triangleright m is the number of voxels
           V_n \leftarrow \text{Voxel } n \text{ containing event points}
 2:
 3:
           t \leftarrow \text{Upper threshold for MST of } V_n
           for \delta = 0 to t do
                                              \triangleright Threshold range from 0 to t
 4:
                 \Phi_{\delta} \leftarrow \text{Node degree distribution at threshold } \delta
 5:
                 \hat{\Phi}_{\delta} \leftarrow \Phi_{\delta} / \max(\Phi_{\delta})
                                                       ▶ Normalize distribution
 6:
                \sigma_{\delta}^2 \leftarrow \text{Variance of } \hat{\Phi}_{\delta}
 7:
 8:
          T_n \leftarrow \arg\max_{\delta}(\sigma_{\delta}^2)
                                                            9:
11: return \{T_n\}_{n=0}^{m-1}
                                        > Set of optimal thresholds for all
```

significant connections, the threshold  $\delta$  is dynamically adjusted based on the normalized degree distribution of the graph. This process involves calculating the degree distribution for each voxel at various values of  $\delta$ , where the distribution captures the number of connections within each voxel and reflects its relevance in the graph structure. This process is illustrated in Algorithm 2

The degree distribution is then normalized to a range of [0,1] to ensure consistency across different scales. This normalization is achieved by dividing the degree distribution  $\Phi_{\delta}$  by its maximum value:

$$\Phi_{\delta} = \frac{\Phi_{\delta}}{\max(\Phi_{\delta})}.\tag{9}$$

After normalization, the variance  $\sigma_{\delta}^2$  of the normalized degree distribution is computed for each threshold  $\delta$ . This variance reflects the spread of the degree distribution, providing insight into the diversity of voxel connections within the graph.

The optimal threshold T is determined by selecting the value of  $\delta$  that maximizes the variance  $\sigma_{\delta}^2$ , formulated as:

$$T = \arg\max_{\xi} (\sigma_{\delta}^2). \tag{10}$$

Selecting the threshold that maximizes variance ensures that T preserves the most meaningful connections between nodes while effectively filtering out noise.

The denoised graph  $G(T) = (\mathcal{V}, \mathcal{E}_T)$  is generated by retaining only edges with weights  $w_{ij} \leq T$ . This filtering process removes irrelevant edges, enhancing the graph's structure. The components of G(T) are:

• Nodes (Event Points): Each node  $i \in \mathcal{V}$  represents an event point with feature vector:

$$f_i = (x_i, y_i, t_i, p_i).$$
 (11)

• Edges (Connections): The edge set  $\mathcal{E}_T$  contains pairs (i, j) satisfying:

$$\mathcal{E}_T = \{(i,j) \mid w_{ij} \le T\}. \tag{12}$$

• Edge Weights: Each edge weight  $w_{ij}$  (see Equation (6))

#### C. AW-GATCN Network Architecture

1) Graph Convolutional Layer with Attention Mechanism: In AW-GATCN, the graph convolutional layer employs an attention mechanism that dynamically adjusts edge weights, emphasizing the most relevant connections. The attention weight  $\alpha_{ij}$  for each edge between nodes i and j is computed based on the features of the target node  $f_i$  and its neighboring node  $f_j$ , using a combination of learnable attention parameters and edge weights. In this framework, smaller edge weights indicate stronger correlations, prompting an inverse adjustment in the attention coefficient  $\alpha_{ij}$  based on edge weight, thereby enhancing the model's ability to capture meaningful relationships between nodes.

$$\alpha_{ij} = \frac{\exp\left(\sigma\left(a^{T}[z_{ij}] \cdot \frac{1}{w_{ij}}\right)\right)}{\sum_{k \in \mathcal{N}(i)} \exp\left(\sigma\left(a^{T}[z_{ik}] \cdot \frac{1}{w_{ik}}\right)\right)},$$
 (13)

where  $\sigma$  is the LeakyReLU activation function, a is a learnable parameter vector that captures the importance of neighboring nodes, and  $z_{ij} = Wf_i \parallel Wf_j$  represents the concatenation of the transformed feature vectors of nodes i and j via the weight matrix W. The edge weight  $w_{ij}$ , calculated through our adaptive weighting method, incorporates  $\frac{1}{w_{ij}}$  to emphasize edges with stronger correlations (i.e., smaller weights), thereby improving feature aggregation. In the denominator,  $\mathcal{N}(i)$  denotes the set of neighboring nodes of i, and k iterates over each neighboring node in  $\mathcal{N}(i)$  for normalization.

The attention coefficients  $\alpha_{ij}$  are then used to aggregate features from neighboring nodes, improving each node's representation by focusing on the connections with the highest correlation. The aggregated feature for node i, denoted by  $f_i'$ , is computed as:

$$f_i' = \sum_{j \in \mathcal{N}(i)} \alpha_{ij} W f_j. \tag{14}$$

By dynamically adjusting each neighbor's influence based on the computed attention weights, this approach allows the network to selectively emphasize informative connections while suppressing less relevant ones. This results in a more accurate representation of the event-based data, ultimately benefiting tasks like object recognition by enhancing robustness against noise.

### IV. PERFORMANCE EVALUATION

# A. Experimental Setup

We evaluate AW-GATCN on four event-based datasets: N-Caltech101, CIFAR10-DVS, MNIST-DVS, and N-CARS. Event data are segmented using adaptive windows based on normalized density and further voxelized via a square root law to balance spatial-temporal granularity. Graphs are constructed within each voxel with edge weights incorporating Euclidean distance, velocity, angular difference, and polarity consistency. Noise is filtered by maximizing the variance of the normalized degree distribution. AW-GATCN employs graph attention for

	C 1 '	· . •	C	1	
IARIH I: Comparison	Of Object	recognition accilrac	v across tour	datacate memo	various methods
TABLE I: Comparison	OI OINCE	. ICCOSIIILIOII aCCUIAC	v across rour	uatasets using	various incurous.

Method	Representation	N-Caltech101	CIFAR10-DVS	MNIST-DVS	N-CARS
H-First	Spike	5.4	7.7	59.5	56.1
Gabor-SNN	Spike	19.6	24.5	82.4	78.9
HOTS	TimeSurface	21.0	27.1	80.3	62.4
HATS	TimeSurface	64.2	52.4	98.4	90.2
DART	TimeSurface	66.4	65.8	98.5	-
YOLE	VoxelGrid	70.2	-	96.1	92.7
AsyncNet	VoxelGrid	74.5	66.3	99.4	94.4
NVS-B	Graph	67.0	60.2	98.6	91.5
NVS-S	Graph	67.0	60.2	98.6	91.5
EvS-B	Graph	76.1	68.0	99.1	93.1
EvS-S	Graph	76.1	68.0	99.1	93.1
AW-GATCN (Ours)	Graph	83.77	76.79	99.3	96.89

feature aggregation, with fixed weight parameters  $\alpha=0.7$ ,  $\beta=0.1$ ,  $\gamma=0.1$ , and  $\delta=0.1$ . All models are trained using PyTorch with 400 epochs under 5-fold cross-validation. Accuracy is used as the evaluation metric.

# B. Comparison with Graph-based Methods for Object Recognition

With optimal weight parameters ( $\alpha=0.7, \beta=0.1, \gamma=0.1, \delta=0.1$ ), we evaluated our AW-GATCN model against state-of-the-art methods on four event-based object recognition benchmarks: N-Caltech101, CIFAR10-DVS, MNIST-DVS, and N-CARS. N-Caltech101, CIFAR10-DVS, and MNIST-DVS are derived from frame-based datasets by displaying moving images on a monitor and recording events with a fixed camera or monitor. N-Caltech101 matches the original Caltech101 in structure, with 8,246 samples across 101 classes. CIFAR10-DVS contains a sixth of the original CIFAR10 dataset, totaling 60,000 samples (6,000 per class). MNIST-DVS uses 10,000 symbols from MNIST, displayed at three scales, for a total of 30,000 samples. In contrast, N-CARS is captured directly with an event camera in real-world scenes, containing 12,336 car and 11,693 non-car samples.

As shown in Table I, AW-GATCN achieved top-tier accuracy across all datasets, significantly outperforming existing approaches. Specifically, our model attained recognition accuracies of 83.77%, 76.79%, 99.3%, and 96.89% on N-Caltech101, CIFAR10-DVS, MNIST-DVS, and N-CARS, respectively. On challenging datasets such as N-Caltech101 and CIFAR10-DVS, AW-GATCN outperformed previous methods by a substantial margin, underscoring its robustness and effectiveness in complex, asynchronous event-based environments, where noise and heterogeneous data often pose significant challenges.

The results validate the effectiveness of the selected weight parameters from Experiment 1, which balanced edge factors to optimize both feature representation and noise reduction. The adaptive weighting approach allows AW-GATCN to dynamically adjust to diverse data characteristics, enabling it to perform well across different data domains.

In summary, these findings demonstrate the efficacy of AW-GATCN asynchronous event processing, establishing it as a robust and accurate model for event-based object recognition. By achieving high accuracy and resilience to noise across diverse settings, AW-GATCN shows significant potential to advance the state-of-the-art in event-driven applications.

#### C. Ablation Study

1) Determination of Weight Parameters: We conducted experiments to optimize the weight parameters for multi-factor edge weighting across three datasets: N-END, D-END [30], and N-CARS, aiming to identify the optimal parameter configuration to maximize classification accuracy. The Event Noisy Dataset (END) consists of two parts: D-END (Daytime) and N-END (Night), providing diverse conditions to evaluate performance across different lighting environments.

The four evaluated weight configurations are as follows:

- Comb 1:  $\alpha = 1, \beta = 0, \gamma = 0, \delta = 0$
- Comb 2:  $\alpha = 0.8, \beta = 0.1, \gamma = 0.05, \delta = 0.05$
- Comb 3:  $\alpha = 0.7, \beta = 0.1, \gamma = 0.1, \delta = 0.1$
- Comb 4:  $\alpha = 0.6, \beta = 0.2, \gamma = 0.1, \delta = 0.1$

These configurations reflect our approach, which leverages multiple factors for noise reduction rather than relying solely on Euclidean distance. Euclidean distance retains a higher weight as the primary criterion for noise determination, while additional factors provide supplementary information. Comb 2 assigns the highest auxiliary weight to velocity vector difference as it is the most relevant secondary factor. Comb 1, using only Euclidean distance, serves as a baseline.

Each configuration was tested using five-fold cross-validation, with the dataset split into five folds. Each fold was used as a test set once, while the remaining four folds formed the training set. The mean accuracy across the five runs was recorded for each configuration.

As shown in Table II, the configuration in Comb 3 achieved the highest accuracy on the D-END and N-CARS datasets, with scores of 93.65% and 96.89%, respectively. This balanced configuration, where Euclidean distance serves as the primary factor and auxiliary factors assist with noise reduction, proved

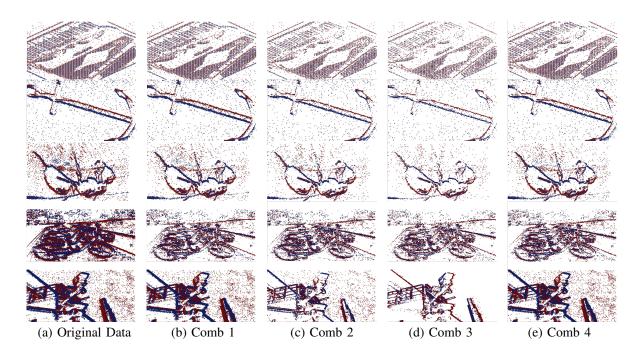


Fig. 3: Compared to Original Data, Comb 3 effectively reduces the number of event points while retaining the primary structure of the recognized object.

TABLE II: Recognition accuracy (%) for various weight parameter combinations across different datasets (D-END, N-CARS, and N-END).

Dataset	Comb 1	Comb 2	Comb 3	Comb 4
D-END N-CARS	89.15 90.63	91.73 94.64	93.65 96.89	90.81 91.57
N-END	79.52	81.71	80.40	80.97

particularly effective for complex datasets. For the N-END dataset, Comb 2 achieved slightly higher accuracy at 81.71%, suggesting that emphasizing Euclidean distance is beneficial for datasets with similar characteristics.

As shown in Figure 3, Comb 3 effectively reduces the number of event points compared to other configurations, achieving an optimal balance between noise reduction and structural preservation. This configuration retains the essential characteristics of the object, allowing the network to capture critical spatial and temporal patterns necessary for accurate recognition. By selectively filtering out redundant or irrelevant data points, Comb 3 creates a more concise and refined representation of events, enhancing recognition performance without compromising the object's primary structure. This balance between noise reduction and structural integrity accounts for Comb 3's superior accuracy across most datasets.

These results indicate that while Comb 3's balanced multifactor approach generally provides optimal performance, adjusting the weight distribution, as in Comb 2, may yield improvements for specific datasets. Overall, a balanced configuration combining Euclidean distance with secondary factors is robust across diverse data conditions.

2) Verification of Denoising Effectiveness: To assess the impact of denoising, we conducted a comparative test of recognition accuracy using Comb 3, evaluating performance with and without denoising across the D-END, N-CARS, and N-END datasets. As shown in Table III, the denoising process significantly enhances recognition accuracy on all datasets. Specifically, denoising improved accuracy by 17.19% on D-END and 19.57% on N-CARS, demonstrating the robustness of AW-GATCN in complex background scenarios where noise can obscure essential features.

For N-END, which primarily consists of data captured at night, the low-light conditions make the structure of objects less distinct, creating a complex noise environment. Denoising led to a 12.2% improvement in accuracy, which, although smaller compared to more challenging datasets, underscores the generalizability of the denoising approach across various data complexities. By reducing noise interference, our model facilitates more accurate feature extraction, allowing the attention mechanism to prioritize meaningful connections over noise, which is essential for object recognition tasks in low-visibility scenarios.

These results validate the role of denoising in the AW-GATCN model, showing that the adaptive noise reduction strategy not only improves robustness in complex environments but also contributes to higher recognition accuracy across diverse data conditions. This ablation study underscores the effectiveness of incorporating denoising within a graph-based neural network framework, emphasizing its impact on model performance.

TABLE III: Recognition accuracy (%) with and without denoising, using the Comb 3 weight configuration on D-END, N-CARS, and N-END datasets.

Comb Group	D-END	N-CARS	N-END
With Denoising	93.65	96.89	80.40
Without Denoising	76.46	77.32	68.2
Improvement	17.19	19.57	12.2

#### V. CONCLUSIONS

We introduced AW-GATCN, an Adaptive Weighted Graph Attention Convolutional Network tailored for event-based data processing, excelling in denoising and object recognition. By integrating adaptive event point segmentation, multifactor edge weighting, and an adaptive graph formulation-based noise reduction approach, AW-GATCN achieves superior accuracy and robustness, especially on noisy, heterogeneous event camera data. Experimental results show that AW-GATCN outperforms state-of-the-art methods with significant accuracy gains on challenging datasets. The optimized weight parameters and attention mechanism effectively prioritize essential connections, capturing spatiotemporal relationships that enhance noise resilience, feature aggregation, and recognition performance.

#### REFERENCES

- R. Berner et al., "A 240×180 10mW 12us latency sparse-output vision sensor for mobile applications," in 2013 Symposium on VLSI Circuits, IEEE, pp. C186–C187, 2013.
- [2] M. Tan and Q. V. Le, "EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks," arXiv, Cornell University, 2019. [Online]. Available: https://arxiv.org/abs/1905.11946. doi: 10.48550/arxiv.1905.11946.
- [3] D. Gehrig et al., "End-to-End Learning of Representations for Asynchronous Event-Based Data," arXiv, Cornell University, 2019. [Online]. Available: https://arxiv.org/abs/1904.08245. doi: 10.48550/arxiv.1904.08245.
- [4] H. Rebecq et al., "Events-to-Video: Bringing Modern Computer Vision to Event Cameras," arXiv, Cornell University, 2019. [Online]. Available: https://arxiv.org/abs/1904.08298. doi: 10.48550/arxiv.1904.08298.
- [5] M. Cannici et al., "A Differentiable Recurrent Surface for Asynchronous Event-Based Data," arXiv, Cornell University, 2020. [Online]. Available: https://arxiv.org/abs/2001.03455. doi: 10.48550/arxiv.2001.03455.
- [6] A. Sironi et al., "HATS: Histograms of Averaged Time Surfaces for Robust Event-based Object Classification," arXiv, Cornell University, 2018. [Online]. Available: https://arxiv.org/abs/1803.07913. doi: 10.48550/arxiv.1803.07913.
- [7] X. Lagorce et al., "HOTS: A Hierarchy of Event-Based Time-Surfaces for Pattern Recognition," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 39, no. 7, pp. 1346–1359, 2017. doi: 10.1109/TPAMI.2016.2574707.
- [8] G. Orchard et al., "HFirst: A Temporal Approach to Object Recognition," 2015. [Online]. Available: https://arxiv.org/abs/1508.01176. doi: 10.48550/arxiv.1508.01176.
- [9] Q. Liu et al., "Effective AER Object Classification Using Segmented Probability-Maximization Learning in Spiking Neural Networks," arXiv, Cornell University, 2020. [Online]. Available: https://arxiv.org/abs/2002.06199. doi: 10.48550/arxiv.2002.06199.
- [10] A. Mitrokhin et al., "Learning Visual Motion Segmentation Using Event Surfaces," in 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), IEEE, pp. 14402–14411, 2020. doi: 10.1109/CVPR42600.2020.01442.
- [11] Y. Bi et al., "Graph-Based Object Classification for Neuromorphic Vision Sensing," in 2019 IEEE/CVF International Conference on Computer Vision (ICCV), IEEE, pp. 491–501, 2019. doi: 10.1109/ICCV.2019.00058.

- [12] Q. Wang, Y. Zhang, J. Yuan, and Y. Lu, "Space-Time Event Clouds for Gesture Recognition: From RGB Cameras to Event Cameras," in 2019 IEEE Winter Conference on Applications of Computer Vision (WACV), Waikoloa, HI, USA, pp. 1826–1835, 2019. doi: 10.1109/WACV.2019.00199. Keywords: Cameras, Three-dimensional displays, Cloud computing, Streaming media, Gesture recognition, Feature extraction, Real-time systems.
- [13] M. Cook et al., "Interacting Maps for Fast Visual Interpretation," in The 2011 International Joint Conference on Neural Networks, IEEE, pp. 770–776, 2011. doi: 10.1109/IJCNN.2011.6033299.
- [14] W. L. Hamilton, R. Ying, and J. Leskovec, "Inductive Representation Learning on Large Graphs," arXiv, Cornell University, 2017. [Online]. Available: https://arxiv.org/abs/1706.02216. doi: 10.48550/arxiv.1706.02216.
- [15] P. Veličković et al., "Graph Attention Networks," arXiv, 2017. [Online]. Available: https://arxiv.org/abs/1710.10903.
- [16] Y. Wang et al., "Event-Stream Representation for Human Gaits Identification Using Deep Neural Networks," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 44, no. 7, pp. 3436–3449, 2022. doi: 10.1109/TPAMI.2021.3054886.
- [17] G. Orchard et al., "A Spiking Neural Network Architecture for Visual Motion Estimation," in 2013 IEEE Biomedical Circuits and Systems Conference (BioCAS), IEEE, pp. 298–301, 2013. doi: 10.1109/Bio-CAS.2013.6679698.
- [18] C. Ye et al., "Unsupervised Learning of Dense Optical Flow, Depth and Egomotion from Sparse Event Data," arXiv, Cornell University, 2018. [Online]. Available: https://arxiv.org/abs/1809.08625. doi: 10.48550/arxiv.1809.08625.
- [19] S. M. M. I et al., "Event-based High Dynamic Range Image and Very High Frame Rate Video Generation using Conditional Generative Adversarial Networks," 2018. [Online]. Available: https://arxiv.org/abs/1811.08230. doi: 10.48550/arxiv.1811.08230.
- [20] N. Messikommer et al., "Event-based Asynchronous Sparse Convolutional Networks," arXiv, Cornell University, 2020. [Online]. Available: https://arxiv.org/abs/2003.09148. doi: 10.48550/arxiv.2003.09148.
- [21] B. Graham, M. Engelcke, and L. van der Maaten, "3D Semantic Segmentation with Submanifold Sparse Convolutional Networks," in 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition, IEEE, pp. 9224–9232, 2018. doi: 10.1109/CVPR.2018.00961.
- [22] J. Xie et al., "Generative PointNet: Deep Energy-Based Learning on Unordered Point Sets for 3D Generation, Reconstruction and Classification," arXiv, Cornell University, 2020. [Online]. Available: https://arxiv.org/abs/2004.01301. doi: 10.48550/arxiv.2004.01301.
- [23] C. R. Qi et al., "PointNet++: Deep Hierarchical Feature Learning on Point Sets in a Metric Space," arXiv, Cornell University, 2017. [Online]. Available: https://arxiv.org/abs/1706.02413. doi: 10.48550/arxiv.1706.02413.
- [24] Y. Sekikawa, K. Hara, and H. Saito, "EventNet: Asynchronous Recursive Event Processing," in 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), IEEE, pp. 3882–3891, 2019. doi: 10.1109/CVPR.2019.00401.
- [25] G. Gallego et al., "Event-Based Vision: A Survey," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 44, no. 1, pp. 154–180, 2022. doi: 10.1109/TPAMI.2020.3008413.
- [26] B. Alwaely and C. Abhayaratne, "AGSF: Adaptive Graph Formulation and Hand-Crafted Graph Spectral Features for Shape Representation," *IEEE Access*, vol. 8, pp. 182260–182272, 2020. doi: 10.1109/AC-CESS.2020.3028696.
- [27] G. Orchard et al., "Converting Static Image Datasets to Spiking Neuromorphic Datasets Using Saccades," arXiv, Cornell University, 2015. [Online]. Available: https://arxiv.org/abs/1507.07629. doi: 10.48550/arxiv.1507.07629.
- [28] H. Li et al., "CIFAR10-DVS: An Event-Stream Dataset for Object Classification," Frontiers in Neuroscience, vol. 11, MAY, pp. 309–309, 2017. doi: 10.3389/fnins.2017.00309.
- [29] T. Serrano-Gotarredona and B. Linares-Barranco, "Poker-DVS and MNIST-DVS. Their History, How They Were Made, and Other Details," *Frontiers in Neuroscience*, vol. 9, pp. 481–481, 2015. doi: 10.3389/fnins.2015.00481.
- [30] S. Ding et al., "E-MLB: Multilevel Benchmark for Event-Based Camera Denoising," *IEEE Transactions on Multimedia*, vol. 26, pp. 1–12, 2024. doi: 10.1109/TMM.2023.3260638.