# **Re-coding for Uncertainties: Edge-awareness Semantic Concordance for Resilient Event-RGB Segmentation**

Nan Bao<sup>1</sup>, Yifan Zhao<sup>1\*</sup>, Lin Zhu<sup>2</sup>, Jia Li<sup>1\*</sup>

<sup>1</sup>State Key Laboratory of Virtual Reality Technology and Systems, SCSE & QRI, Beihang University <sup>2</sup>School of Computer Science and Technology, Beijing Institute of Technology {nbao, zhaoyf, jiali}@buaa.edu.cn, linzhu@bit.edu.cn

# **Abstract**

Semantic segmentation has achieved great success in ideal conditions. However, when facing extreme conditions (e.g., insufficient light, fierce camera motion), most existing methods suffer from significant information loss of RGB, severely damaging segmentation results. Several researches exploit the high-speed and high-dynamic event modality as a complement, but event and RGB are naturally heterogeneous, which leads to feature-level mismatch and inferior optimization of existing multi-modality methods. Different from these researches, we delve into the edge secret of both modalities for resilient fusion and propose a novel Edge-awareness Semantic Concordance framework to unify the multi-modality heterogeneous features with latent edge cues. In this framework, we first propose Edge-awareness Latent Re-coding, which obtains uncertainty indicators while realigning event-RGB features into unified semantic space guided by re-coded distribution, and transfers event-RGB distributions into re-coded features by utilizing a pre-established edge dictionary as clues. We then propose Re-coded Consolidation and Uncertainty Optimization, which utilize re-coded edge features and uncertainty indicators to solve the heterogeneous event-RGB fusion issues under extreme conditions. We establish two synthetic and one real-world event-RGB semantic segmentation datasets for extreme scenario comparisons. Experimental results show that our method outperforms the state-of-the-art by a 2.55% mIoU on our proposed DERS-XS, and possesses superior resilience under spatial occlusion. Our code and datasets are publicly available at https://github.com/iCVTEAM/ESC.

# 1 Introduction

Being widely used in autonomous driving, medical imaging, geospatial analysis, and industrial inspection, semantic segmentation aims to resolve the semantics of visual objects, assigning category labels to each pixel in the image [8]. When facing extreme conditions due to diversity and complexity in the wild, conventional single-RGB semantic segmentation faces challenges of corrupted results, suffering from significant information loss. This has led to the exploration of leveraging information from multiple modalities for semantic segmentation [46, 50, 51].

We investigate the problem of leveraging event and RGB for semantic segmentation under extreme conditions, focusing on inferior optimization issues in modality imbalance and failure situations. Dynamic vision sensor [24, 6, 38], commonly known as event camera, responds to brightness changes and generates events for each pixel asynchronously and independently. This unique mechanism gives it many advantages beyond conventional cameras, such as high dynamic range, high temporal resolution, low latency, and low power consumption [10]. Therefore, event data are widely used

<sup>\*</sup>Correspondence should be addressed to Yifan Zhao and Jia Li. Website: https://cvteam.buaa.edu.cn/

in tasks that are hard to solve with conventional images alone, such as HDR image reconstruction [31, 35, 53], motion deblurring [29, 17, 32], and low-light enhancement [16, 23, 59, 58]. As shown in fig. 1, image suffers from severe information loss under extreme conditions due to a low signal-to-noise ratio, while events clearly show the motion edge of vehicles. It becomes feasible to complement the lost information of RGB modality by utilizing event modality.

However, existing event-RGB semantic segmentation methods [50, 51, 46] do not consider the heterogeneous properties of event and RGB modality. Therefore, although the naive fusion strategy has achieved some improvements, it is difficult to handle featurelevel mismatch and inferior optimization issues, especially in modality imbalance and failure situations. To overcome the above problems of heterogeneous event and RGB, we find semantic edge as an intermediate commonality for both.

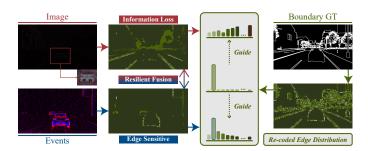


Figure 1: Edge-awareness Semantic Concordance for event-RGB fusion. RGB suffers from severe information loss under extreme conditions, while events are sensitive to edges in motion, complementing the lost information. Heterogeneous properties of event and RGB lead to feature-level mismatch and inferior optimization of existing methods. Our ESC framework utilizes semantic edge as an intermediate commonality for a more resilient fusion.

Existing studies [22, 56, 45] have proven that edge-awareness is beneficial for RGB segmentation. Intuitively, events highlight edges in motion, and RGB gradients reveal edge cues. Through statistics in section 3.1, we indeed find a strong correlation between events and semantic edge. **Semantic edge serves as a bridge, guiding the heterogeneous event and RGB to embed into a unified semantic space.** By utilizing information of semantic edge as bridge, we successfully realign the heterogeneous event and RGB into the unified semantic space to jointly optimize their edge information, while consolidating the image contextual information with semantic edge information as crucial clue.

In this paper, we propose Edge-awareness Semantic Concordance (ESC), a novel multi-modality learning framework for event-RGB semantic segmentation. ESC utilizes a shared discrete embedding space, creating an edge dictionary containing basic semantic elements from semantic edge. We introduce Edge-awareness Latent Re-coding for discrete latent space modeling and transferring bidirectionally, namely re-coding. The re-coded edge features are utilized for information consolidation, and the re-coded edge distribution enables unified realignment through cross-entropy supervision. Uncertainty indicators are derived from modality distributions for joint optimization. Re-coded Consolidation and Uncertainty Optimization are designed to achieve the above processes for resilient fusion. Prior work [46] assesses event-RGB segmentation using RGB-pseudo-labeled datasets (e.g., DDD17 [4], DSEC-Semantic[11, 36]), leading to potentially unreliable results. To address this, we introduce synthetic DERS-XS and real-world DERS-XR, featuring low-light RGB, noisy events, and true-labels for extreme scenario comparisons. We further adapt DSEC-Semantic into an extreme variant, DSEC-Xtrm, to mitigate direct dependence of pseudo-labels on original RGB. Experiments on above datasets show that our method achieves better performance and is more resilient under extreme conditions compared to existing multi-modality methods. To the best of our knowledge, this is the first work to assess model resilience via spatial occlusion evaluation without any fine-tuning.

The contributions of our work are summarized as follows:

- 1) We propose Edge-awareness Semantic Concordance (ESC), a novel multi-modality framework that exploits supervision over re-coded distribution to realign heterogeneous event and RGB into unified semantic space, jointly optimizing them based on uncertainties derived from modality distributions.
- 2) We propose three modules, namely Edge-awareness Latent Re-coding (ELR), Re-coded Consolidation (RC), and Uncertainty Optimization (UO). ELR re-codes features and distributions bidirectionally, while RC and UO utilize the re-coded features and uncertainties for a resilient fusion.
- 3) We establish two synthetic and one real-world event-RGB semantic segmentation datasets for extreme scenario comparisons. Experimental results show that our method outperforms the state-of-the-art methods and possesses superior resilience under extreme conditions including occlusion.

# 2 Related work

# 2.1 Event-based semantic segmentation

Event data has recently been applied to semantic segmentation tasks. Ev-SegNet [1] introduces semantic segmentation for event data by proposing a 6-channel image-like representation and applying CNN architectures on the DDD17 dataset [4]. EvDistill [40] trains a student network on unlabeled event data via knowledge distillation from a large image-based teacher network. CMDA [43] proposes an unsupervised domain adaption segmentation framework to transfer daytime RGB knowledge to nighttime event domain. ESS [36] leverages labeled images for training on unlabeled event data through unsupervised domain adaptation. EvSegFormer [15] introduces a posterior attention module to utilize prior knowledge from event data, and HPL-ESS [18] proposes a hybrid pseudo-labeling framework to mitigate noisy labels in unsupervised event-based segmentation. ESEG [55] is a uni-modality event-based segmentation framework that exploits edge semantics to provide explicit guidance toward the regions of interest. ISSAFE [52] leverages events to assist segmentation under accident scenes, and HALSIE [5] features a hybrid dual-encoder scheme with SNN and ANN for efficient segmentation. Recent works demonstrate the feasibility of leveraging event data, while most works do not fully explore its unique characteristics, limiting its advantages over conventional RGB.

#### 2.2 Event-assisted vision tasks

Event data can be utilized for assisting conventional vision tasks due to its high-speed and high-dynamic capacity. Pan *et al.* [29] propose an event-based double integral model to restore sharp video from a single blurry frame with events. Jiang *et al.* [17] recover sharp videos with events based on a recurrent neural network. Shang *et al.* [32] utilize events for non-consecutively frames deblurring. Liang *et al.* [23] and Liu *et al.* [26] leverage event data to guide low-light video enhancement. Jiang *et al.* [16] propose a joint framework to reconstruct clear images from underexposed frames and event streams. Shi *et al.* [33] utilize paired images and event streams to estimate monocular depth under night conditions. Li *et al.* [21] propose an event-based low-light video object segmentation framework. Qi *et al.* [30] introduce events into neural radiance fields for novel view sharp image rendering. Geng *et al.* [12] introduce events into visible and infrared fusion task. There are also several cross-modality contrastive pretraining approaches, such as Yang *et al.* [48], Yao *et al.* [49], and Wu *et al.* [42], aiming to acquire informative and effective pretrained backbones for both event and RGB. Prevailing research proves the effectiveness of event-assisted tasks, while the inferior optimization issue of heterogeneous event and RGB under extreme conditions remains unexplored.

# 2.3 Inter-modality Fusion

Inter-modality fusion is the core issue of multi-modality tasks. How to obtain better-fused features has become an enduring research topic. Zhang et al. [50, 51] aim to achieve a general cross-modality segmentation model for arbitrary modalities, including event modality. Xie et al. [46] propose a modality recalibration and fusion module to recalibrate and then aggregate events and image features at multiple stages. Other works only focus on fusion techniques without specifying a specific task. Wang et al. [41] detect tokens with less information dynamically and substitute them with aggregated features projected from another modality. Jia et al. [14] introduces noise embeddings into proposed inter-modality attention module to improve interaction between features of multi-modality pixel-wise. Zhao et al. [54] utilized extra edge cues for event-RGB stereo. Several approaches have also emerged, including Kim et al. [19] and Lang et al. [20], to address the problem of incomplete modality inputs. In vision-language models, discrete representation learning with shared embedding space is becoming popular. Liu et al. [25] propose a representation learning paradigm that contains a discretized embedding space shared across two different modalities such as video and audio. Xia et al. [44] propose a framework that obtains mutual semantic information from different modalities by modality feature reconstruction. Zheng et al. [57] propose an iterative learning paradigm for tuning large language models into multi-modality LLM. Zhou et al. [60] draw on the concept of shared latent space and first introduce it into domain adaptation vision task of nighttime optical flow estimation.

Inspired by the above works, we propose an Edge-awareness Semantic Concordance framework to model a shared discrete latent edge space and optimize events and image features into the unified semantic space based on the re-coded edge distribution. By edge-awareness latent re-coding, we obtain re-coded edge features and uncertainties, which are utilized for inter-modality resilient optimization.

# 3 Method

Naive fusion strategy fails to integrate heterogeneous event and RGB under extreme conditions. We propose an Edge-awareness Semantic Concordance framework to address this. To prove the rationality, we first analyze event edge characteristics in section 3.1. We then establish an edge dictionary as a preliminary in section 3.2. Based on this dictionary, Edge-awareness Latent Re-coding (section 3.3) transforms edge distribution and features bi-directionally, namely re-coding. Re-coded edge distribution is utilized for feature-level unified realignment through supervision. Re-coded Consolidation and Uncertainty Optimization (section 3.4 and section 3.5) utilizes re-coded edge features and uncertainties derived from modality distributions for a resilient fusion. Since labels from DSEC-Semantic are unreliable, we construct three new datasets for reliable evaluation in section 3.6.

## 3.1 Edge characteristic of events

Event camera is a bio-inspired sensor that triggers event signals asynchronously when light intensity changes at each pixel. Specifically, as eq. (1) shows, an event  $\mathbf{e} = \langle \mathbf{x}, t, p_{\mathbf{x},t} \rangle$  is triggered when pixel  $\mathbf{x} = \langle x, y \rangle$  perceives a change in light intensity I that reaches threshold  $\Theta$  in the logarithmic domain at time t, where  $p_{\mathbf{x},t}$  means polarity of light intensity change in logarithm domain. Triggered events from  $t_{start}$  to  $t_{end}$  form an event stream  $\{\mathbf{e}_i = \langle \mathbf{x}_i, t_i, p_i \rangle\}_{t_{start} < t_i \leq t_{end}}$ .

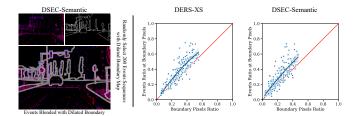


Figure 2: Correlation between events and semantic edge. We randomly select 200 event sequences with dilated boundary map from true-labeled DERS-XS and real-world DSEC-Semantic, counting the ratio of edge pixels to all pixels and the ratio of events at edge pixels to all events, respectively. For both datasets, as the area of edge expands, the events ratio is always greater than the boundary ratio. This exhibits a strong correlation between events and semantic edge under different conditions.

$$p_{\mathbf{x},t} = \begin{cases} +1, & \log(I_{\mathbf{x},t}) - \log(I_{\mathbf{x},t-\Delta t}) > \Theta, \\ -1, & \log(I_{\mathbf{x},t}) - \log(I_{\mathbf{x},t-\Delta t}) < -\Theta. \end{cases}$$
(1)

We demonstrate the correlation between events and semantic edge (*i.e.*, segmentation boundary) through statistics in fig. 2. We randomly select 200 event sequences with dilated boundary map from true-labeled DERS-XS and real-world DSEC-Semantic, counting the ratio of edge pixels to all pixels of the whole plane and the ratio of events falling on edge pixels to all events of the whole sequence. Statistical results show that as the area of edge pixels expands, the events ratio is always greater than the boundary ratio for both datasets. This demonstrates that events tend to cluster at areas of semantic edge under different conditions, exhibiting a strong correlation between events and semantic edge, which supports our utilization of semantic edge as a bridge for heterogeneous event and RGB. Details of statistical process of event-edge correlation with more analyses can be seen in appendix C.

## 3.2 Edge dictionary as intermediate semantic clues across modalities

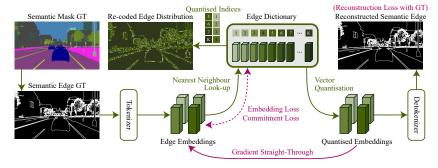


Figure 3: **Establishment of edge dictionary.** We establish our edge dictionary based on a VQ-VAE architecture. Semantic edge is retrieved from the semantic mask ground truth and leveraged for learning its discrete latent representations as an edge dictionary, which serves as intermediate clues across heterogeneous event and RGB.

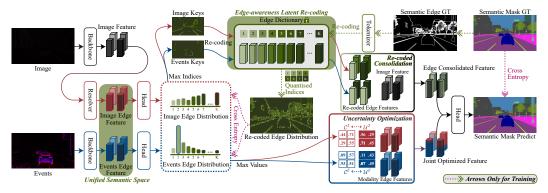


Figure 4: The overall architecture of our Edge-awareness Semantic Concordance (ESC). ESC contains a pre-established edge dictionary and three key modules, namely Edge-awareness Latent Re-coding (ELR), Re-coded Consolidation (RC), and Uncertainty Optimization (UO). Based on the pre-trained edge dictionary, ELR transfers edge embeddings into re-coded distribution and modality distribution into re-coded features. RC consolidates edge information with re-coded features. UO jointly optimizes modality edge features with uncertainties. Features from RC and UO are concatenated for final semantic mask prediction.

To utilize semantic edge as intermediate clues, we establish an edge dictionary, which is a discrete latent embedding space derived from semantic edge, containing basic semantic elements of edge and shared by heterogeneous event and RGB. The establishment of our edge dictionary is shown in fig. 3 and based on a VQ-VAE [39] architecture, which is originally used for representation learning and used by [25, 28, 44, 57] to model shared discrete latent space from inputs.

We first retrieve semantic edge from semantic mask ground-truth by a mean filter followed by an indicator function. Given a semantic mask  $\mathbf{M} \in \{1, \cdots, c\}^{H \times W}$ , boundary map  $\mathbf{B} \in \{0, 1\}^{H \times W}$  can be obtained by  $\mathbf{B} = \mathbb{I}_{\mathbf{M} \neq \text{Mean-Filter}(\mathbf{M})}(\mathbf{M})$ , where c is number of categories in semantic mask, and  $\mathbb{I}$  is indicator function.

We define edge dictionary as  $\Delta = \{\langle k, v(k) \rangle | k \in \{1, \cdots, K\} \}$ , where K is the number of items (i.e. quantised vectors) that edge dictionary contains, and  $v(\cdot)$  is the query function as  $v(k) \in \mathbb{R}^n$  selects the k-th item of edge dictionary with n-dimension. As fig. 3 shows, the tokenizer  $f_T$  takes boundary map  $\mathbf{B}$  as input, producing edge embeddings  $\mathbf{\Gamma} = f_T(\mathbf{B}) \in \mathbb{R}^{H' \times W' \times n}$ , which have downsampled spatial size  $H' \times W'$  and n channels. Items in edge dictionary are selected by nearest neighbour look-up method as  $\mathbf{\Gamma}' = v(\hat{\mathbf{K}}) = v(\arg\min_{\mathbf{K} \in \{1, \cdots, K\}^{H' \times W'}} \|\mathbf{\Gamma} - v(\mathbf{K})\|_2^2)$ , where  $\hat{\mathbf{K}}$  contains the queried indices and  $\mathbf{\Gamma}'$  is the quantised edge embeddings. Reconstructed boundary map  $\mathbf{B}'$  is obtained by  $\mathbf{B}' = f_{T'}(\mathbf{\Gamma}')$ , where  $f_{T'}$  is the detokenizer.

To ensure the edge dictionary contains all basic information of semantic edge, we need to ensure the boundary map is reconstructed flawlessly while items in edge dictionary are close enough to edge embeddings. Thus, we adopt the training objective with reconstruction loss, embedding loss, and commitment loss as  $L_{dict} = \|\mathbf{B} - \mathbf{B}'\|_2^2 + \|v(\hat{\mathbf{K}}) - \mathrm{sg}(\mathbf{\Gamma})\|_2^2 + \alpha\|\mathrm{sg}(v(\hat{\mathbf{K}})) - \mathbf{\Gamma}\|_2^2$ , where  $\mathrm{sg}(\cdot)$  means stop gradient, and  $\alpha$  is a constant of commitment loss weight. To make the reconstruction loss propagate back to tokenizer, a gradient straight-through technique is adopted, which directly assigns the gradient from  $\mathbf{\Gamma}'$  to  $\mathbf{\Gamma}$ . Details of edge dictionary training process can be seen in appendix  $\mathbf{D}$ .

### 3.3 Cross-modality realignment of edge representations via latent re-coding

Re-coding is a key process in our framework, realigning edge representations of heterogeneous event and RGB through re-coded distribution, while also producing re-coded edge features for consolidation. Based on the pre-established edge dictionary, our proposed Edge-awareness Latent Re-coding module transfers edge embeddings into re-coded edge categorical prior distribution and modality posterior distribution into re-coded edge features. This section will discuss the latent re-coding operation of two directions mentioned above, and introduce the optimization objective at the end of this section.

**Re-coding for edge categorical prior distribution.** Given the pre-trained tokenizer  $f_T$  and the pre-established edge dictionary  $\Delta$ , we can re-code any semantic edge **B** to an edge categorical prior

distribution  $q(\mathcal{K}|\mathbf{B})$  as one-hot as follows:

$$q(\mathcal{K}|\mathbf{B}) = \underset{\mathcal{K} \in \{\mathbf{b}_k | k \in \{1, \dots, K\}\}^{H' \times W'}}{\operatorname{arg min}} \|f_T(\mathbf{B}) - v(k)\|_2^2, \tag{2}$$

where  $\mathbf{b}_k$  is the K-dim basis vector with 1 at k-th place.

**Re-coding for edge features.** We first extract features from inputs. Given an image  $\mathcal{I} \in \mathbb{R}^{H \times W \times 3}$  and its corresponding event voxel grid [61]  $\mathcal{E} \in \mathbb{R}^{H \times W \times B}$ , multi-scaled image feature and events feature can be obtained from  $\mathbf{F}^{\mathcal{I}} = f_{\mathcal{I}}(\mathcal{I}), \mathbf{F}^{\mathcal{E}} = f_{\mathcal{E}}(\mathcal{E})$ , where B is number of voxel grid bins,  $\mathbf{F}^{\mathcal{I}}$  and  $\mathbf{F}^{\mathcal{E}}$  denotes the backbones. Image edge features are resolved as  $\mathbf{E}^{\mathcal{I}} = f_{R}(\mathbf{F}^{\mathcal{I}})$ , where  $f_{R}$  is an decouple module adopted from [22] as our edge resolver. As event data naturally highlights the edge information, we keep event features directly as events edge features  $\mathbf{E}^{\mathcal{E}} = \mathbf{F}^{\mathcal{E}}$  without additional processing. Then two edge encoders with the same structure are applied respectively to both  $\mathbf{E}^{\mathcal{I}}$  and  $\mathbf{E}^{\mathcal{E}}$ , in order to encode modality edge features into the same unified semantic space. Two MLP-based classification heads are attached after edge encoders for each modality to predict its modality-specific edge categorical probability distribution  $p(\mathcal{K}|\mathcal{I})$  and  $p(\mathcal{K}|\mathcal{E})$ . This categorical probability distribution indicates the probability of edge dictionary item number K-ary classification at each spatial position.

Given probability distributions  $p(\mathcal{K}|\mathcal{I})$  and  $p(\mathcal{K}|\mathcal{E})$ , we can retrieve image key map  $\mathbf{K}^{\mathcal{I}}$  and events key map  $\mathbf{K}^{\mathcal{E}}$  by

$$\mathbf{K}^{\mathcal{M}} = \underset{k \in \{1, \dots, K\}}{\operatorname{arg max}} p(\mathcal{K} = k | \mathcal{M}), \quad \mathcal{M} \in \{\mathcal{I}, \mathcal{E}\},$$
(3)

where  $\mathbf{K}^{\mathcal{M}} \in \{1,\cdots,K\}^{H' \times W'}$  means the key map of modality  $\mathcal{M}$ , which can either be image modality  $\mathcal{I}$  or events modality  $\mathcal{E}$ . By this step, we select the indices of the maximum probability values at each position of the latent space as key map, which can be utilized to query edge dictionary for obtaining re-coded edge features of the specific modality. The image re-coded edge feature  $\mathbf{\Gamma}^{\mathcal{I}}$  and events re-coded edge feature  $\mathbf{\Gamma}^{\mathcal{E}}$  are obtained by

$$\mathbf{\Gamma}^{\mathcal{M}} = v(\mathbf{K}^{\mathcal{M}}), \quad \mathcal{M} \in \{\mathcal{I}, \mathcal{E}\},$$
(4)

where  $\mathbf{\Gamma}^{\mathcal{M}} \in \mathbb{R}^{H' \times W' \times n}$  are modality-specific edge embeddings queried by the specific key map.

How to optimize ELR and what are the benefits? We optimize our Edge-awareness Latent Recoding module by an objective function based on cross-entropy, which narrows the gap between the edge categorical distribution  $q(\mathcal{K}|\mathbf{B})$  with modality-specific edge categorical probability distribution  $p(\mathcal{K}|\mathcal{I})$  and  $p(\mathcal{K}|\mathcal{E})$  as  $L_{edge} = -\sum q(\mathcal{K}|\mathbf{B})\log(p(\mathcal{K}|\mathcal{I})p(\mathcal{K}|\mathcal{E}))$ , a summation of two cross-entropies. By minimizing this objective function, we can bridge the modality gap and realign the image edge feature  $\mathbf{E}^{\mathcal{I}}$  with events edge feature  $\mathbf{E}^{\mathcal{E}}$  into the same unified semantic space, and make sure the re-coded features  $\mathbf{\Gamma}^{\mathcal{M}}$  represent the edge information of events and image correctly.

# 3.4 Re-coded features for edge information consolidation

Image feature mainly focuses on contextual information and lack of understanding of edge information. Recoded features are utilized for edge information consolidation by the Recoded Consolidation (RC) module. As shown in fig. 5 on the left, RC takes image feature  $\mathbf{F}^{\mathcal{I}}$ , image and events re-coded edge feature  $\mathbf{\Gamma}^{\mathcal{I}}$  and  $\mathbf{\Gamma}^{\mathcal{E}}$  as input, outputs a refined feature namely Edge Consolidated Feature  $\mathbf{\Phi}$ .

In RC, we define two learnable noise embeddings  $\mathbf{N}_K \in \mathbb{R}^n$  and  $\mathbf{N}_V \in \mathbb{R}^n$  to improve fitting ability and enhance learning stability. For image feature  $\mathbf{F}^{\mathcal{I}}$  and re-coded edge features

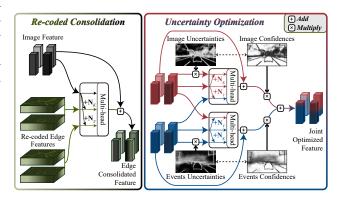


Figure 5: **RC** and **UO**. The two modules utilize an attention-based structure with learnable noise embeddings for a resilient fusion.

 $\Gamma^{\mathcal{I}}$  and  $\Gamma^{\mathcal{E}}$ , RC applies multi-head attention operation on vectors at spatial position  $\langle h, w \rangle$ , calculates and outputs the consolidated features  $\Phi$  as

$$\Phi_{h,w} = [\phi_1, \phi_2, \cdots, \phi_m] \cdot W_O + \mathbf{F}_{h,w}^{\mathcal{I}}, 
\phi_i = \text{Softmax}(Q_i K_i^{\mathcal{I}} / \sqrt{d_k}) \cdot V_i,$$
(5)

where  $Q_i = \mathbf{F}_{h,w}^{\mathcal{I}} \cdot W_{Q_i}$ ,  $K_i = [\mathbf{F}_{h,w}^{\mathcal{I}} + \mathbf{N}_K, \mathbf{\Gamma}_{h,w}^{\mathcal{I}}, \mathbf{\Gamma}_{h,w}^{\mathcal{E}}] \cdot W_{K_i}$ ,  $V_i = [\mathbf{F}_{h,w}^{\mathcal{I}} + \mathbf{N}_V, \mathbf{\Gamma}_{h,w}^{\mathcal{I}}, \mathbf{\Gamma}_{h,w}^{\mathcal{E}}] \cdot W_{V_i}$ . The introduction of noise embedding is inspired by [14], and we develop the following theoretical explanation. In the absence of noise, the query  $(Q_i)$  tends to attend excessively to its own features in the key  $(K_i)$ , thereby suppressing signals from the other source and impeding effective fusion. Introducing noise embeddings mitigates this issue by perturbing the attention space in a controlled, learnable manner, encouraging richer and more balanced cross-modality interactions. The main idea of eq. (5) is to consolidate image feature  $\mathbf{F}^{\mathcal{I}}$  with image re-coded edge feature  $\mathbf{\Gamma}^{\mathcal{I}}$  and events re-coded edge feature  $\mathbf{\Gamma}^{\mathcal{E}}$  by querying  $\mathbf{\Gamma}^{\mathcal{I}}$  and  $\mathbf{\Gamma}^{\mathcal{E}}$  with  $\mathbf{F}^{\mathcal{I}}$  to obtain attention map. The map represents the relevance between  $\mathbf{F}^{\mathcal{I}}$  and  $\mathbf{\Gamma}^{\mathcal{E}}$ , which decides the amount of edge information consolidated by  $\mathbf{F}^{\mathcal{I}}$ . Refined vectors  $\mathbf{\Phi}_{h,w}$  constitute together in accordance with their positions as  $\mathbf{\Phi} \in \mathbf{R}^{H_d \times W_d \times n}$ .

### 3.5 Edge-aware uncertainties for joint optimization

Probability values in edge distribution indicate the confident and uncertain areas of image and events. We leverage this confidence and uncertainty information from edge distribution of image and events for a resilient fusion by the Uncertainty Optimization (UO) module.

Given edge categorical probability distributions  $p(\mathcal{K}|\mathcal{I})$  and  $p(\mathcal{K}|\mathcal{E})$ , we can retrieve confidences and uncertainties of image and events by

$$C^{\mathcal{M}} = \max_{k \in \{1, \dots, K\}} p(\mathcal{K} = k | \mathcal{M}),$$

$$U^{\mathcal{M}} = 1 - C^{\mathcal{M}}, \quad \mathcal{M} \in \{\mathcal{I}, \mathcal{E}\},$$
(6)

where  $\mathcal{C}^{\mathcal{M}}, \mathcal{U}^{\mathcal{M}} \in [0, 1]^{H' \times W'}$  denote the confidences and uncertainties of modality  $\mathcal{M}$ , which can either be image modality  $\mathcal{I}$  or events modality  $\mathcal{E}$ .

Confidences and Uncertainties represent spatial reliability of specific modality, which are utilized as indicators in UO. As shown in fig. 5 on the right, UO takes image edge feature  $\mathbf{E}^{\mathcal{E}}$  and events edge feature  $\mathbf{E}^{\mathcal{E}}$  as input, confidences and uncertainties as indicators, and outputs a refined feature namely Joint Optimized Feature  $\Psi$ .

In UO, we define four learnable noise embeddings  $\mathbf{N}_K^{\mathcal{I}}, \mathbf{N}_K^{\mathcal{E}}, \mathbf{N}_V^{\mathcal{I}}, \mathbf{N}_V^{\mathcal{E}} \in \mathbb{R}^n$  to enhance fitting capability and improve learning robustness. For feature  $\mathbf{E}^{\mathcal{I}}$  and  $\mathbf{E}^{\mathcal{E}}$ , their confidences and uncertainties are  $\mathcal{C}^{\mathcal{I}}, \mathcal{U}^{\mathcal{I}}$  and  $\mathcal{C}^{\mathcal{E}}, \mathcal{U}^{\mathcal{E}}$  respectively. UO applies multi-head attention operation on vectors at spatial position  $\langle h, w \rangle$ , calculates and outputs the optimized feature  $\mathbf{\Psi}$  as

$$\Psi_{h,w} = \frac{C_{h,w}^{\mathcal{I}} \cdot \Psi_{h,w}^{\mathcal{I}}}{C_{h,w}^{\mathcal{I}} + C_{h,w}^{\mathcal{E}}} + \frac{C_{h,w}^{\mathcal{E}} \cdot \Psi_{h,w}^{\mathcal{E}}}{C_{h,w}^{\mathcal{I}} + C_{h,w}^{\mathcal{E}}},$$

$$\Psi_{h,w}^{\mathcal{M}} = [\psi_{1}^{\mathcal{M}}, \psi_{2}^{\mathcal{M}}, \cdots, \psi_{m}^{\mathcal{M}}] \cdot W_{O}^{\mathcal{M}} + \mathbf{E}_{h,w}^{\mathcal{M}},$$

$$\psi_{i}^{\mathcal{M}} = \operatorname{Softmax}(Q_{i}^{\mathcal{M}}(K_{i}^{\mathcal{M}})^{\mathsf{T}}/\sqrt{d_{k}}) \cdot V_{i}^{\mathcal{M}},$$
(7)

where  $Q_i^{\mathcal{M}} = \mathbf{E}_{h,w}^{\mathcal{M}} \cdot W_{Q_i}^{\mathcal{M}}$ ,  $K_i^{\mathcal{M}} = [\mathbf{E}_{h,w}^{\mathcal{M}} + \mathbf{N}_K^{\mathcal{M}}, \mathcal{U}_{h,w}^{\mathcal{M}} \cdot \mathbf{E}_{h,w}^{\mathcal{M}}] \cdot W_{K_i}^{\mathcal{M}}$ ,  $V_i^{\mathcal{M}} = [\mathbf{E}_{h,w}^{\mathcal{M}} + \mathbf{N}_V^{\mathcal{M}}, \mathbf{E}_{h,w}^{\mathcal{M}}] \cdot W_{V_i}^{\mathcal{M}}$ ,  $V_i^{\mathcal{M}} = [\mathbf{E}_{h,w}^{\mathcal{M}} + \mathbf{N}_V^{\mathcal{M}}, \mathbf{E}_{h,w}^{\mathcal{M}}] \cdot W_{V_i}^{\mathcal{M}}$ ,  $V_i^{\mathcal{M}} = [\mathbf{E}_{h,w}^{\mathcal{M}} + \mathbf{N}_V^{\mathcal{M}}, \mathbf{E}_{h,w}^{\mathcal{M}}] \cdot W_{V_i}^{\mathcal{M}}$ ,  $V_i^{\mathcal{M}} = [\mathbf{E}_{h,w}^{\mathcal{M}} + \mathbf{N}_V^{\mathcal{M}}, \mathbf{E}_{h,w}^{\mathcal{M}}] \cdot W_{V_i}^{\mathcal{M}}$ , which decides the optimized image edge feature  $\mathbf{E}^{\mathcal{E}}$  and events edge feature  $\mathbf{E}^{\mathcal{E}}$  based on their confidences at each spatial position. Multiplied by the uncertainty value  $\mathcal{U}^{\mathcal{M}}$ , modality-specific feature exposes its uncertainty to attention map. The map represents the self-uncertainty of  $\mathbf{E}^{\mathcal{I}}$  and  $\mathbf{E}^{\mathcal{E}}$ , which decides the amount of complementary information absorbed from the counter modality. The final feature vector is calculated by normalized confidence weighted summation of inter-modalities feature vectors. Optimized vectors  $\mathbf{\Psi}_{h,w}$  constitute together in accordance with their positions as  $\mathbf{\Psi} \in \mathbf{R}^{H_d \times W_d \times n}$ .

How to optimize our method? Edge consolidate feature  $\Phi$  and joint optimized feature  $\Psi$  are concatenated and input into an MLP-based classification head for semantic mask prediction. Crossentropy is utilized for the supervision of semantic mask prediction as  $L_{pred}$ . The final optimization objective function for our method is  $L = L_{pred} + \beta \cdot L_{edge}$ , where  $\beta$  is a constant of edge loss weight.

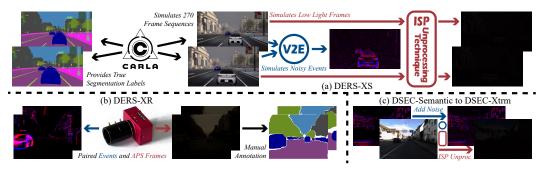


Figure 6: Construction of datasets DERS-XS, DERS-XR, and DSEC-Xtrm for reliable evaluation.

# 3.6 Constructing datasets for reliable evaluation of event-RGB segmentation

Labels from DSEC-Semantic are pseudo-labels directly derived from RGB via [37]. They are useful for event-only tasks, but not reliable for evaluating event-RGB tasks. Using these labels as ground truth implicitly presupposes that the optimal result of event-RGB segmentation is obtained by an RGB-only model, which undermines the unique advantages of events. To address the problem, we construct datasets as below for reliable evaluation. The construction pipelines of datasets DERS-XS, DERS-XR, and DSEC-Xtrm are shown in fig. 6. Details of the datasets can be seen in appendix A.

**DERS-XS**. Dataset of Event-RGB semantic Segmentation under eXtreme conditions Synthetic, abbreviated as DERS-XS, is a true-labeled synthetic event-RGB extreme condition semantic segmentation dataset. **CARLA** [9] **provides true segmentation labels**, and we first use CARLA simulator to simulate 270 frame sequences with segmentation labels of 23 categories, each with 1200 frames, and the size of each frame is  $640 \times 360$ . We then simulate noisy events from the frame sequences using v2e simulator [13] with a shot noise parameter of 5.0 Hz. Low-light frames are simulated by attenuating optical signals and adding shot noise in the RAW domain obtained from the ISP unprocessing technique in [7]. Because differences between adjacent frames are small, which is not conducive to data diversity, we sample data at intervals of 100 frames while discarding other frames. We divide 168 sequences as training set, 12 sequences as validation set, and 90 sequences as test set.

**DERS-XR**. Dataset of Event-RGB semantic Segmentation under eXtreme conditions Real-world, abbreviated as DERS-XR, is a **manually annotated** real-world event-RGB extreme condition semantic segmentation dataset. We use a DAVIS346 [38] to capture paired APS frames and events under extreme lighting conditions, and manually annotate a subset of 240 frames. Of these, 120 frames are randomly selected for fine-tuning, while the remaining 120 frames are used for testing.

**DSEC-Xtrm**. DSEC-Xtrm is an extreme condition semantic segmentation dataset synthesized based on DSEC-Semantic [11, 36]. **To make use of pseudo-labels from DSEC-Semantic while mitigating their direct dependence on RGB**, we apply degradation to the RGB frames. We apply the same low-light image simulation method as DERS-XS to frames and use v2e simulator [13] to generate pure shot noise and add it to events. The degraded frames and events together constitute DSEC-Xtrm.

# 4 Experiments

# 4.1 Implement details

The code is implemented by PyTorch. We first train edge dictionary as a separate stage to obtain pre-trained weights of tokenizer and edge dictionary. We utilize pre-trained MiT-B2 backbone and MiT-B1 backbone of SegFormer [47] for RGB modality and event modality, respectively. The number of categories c is 11. We set the number of items K in edge dictionary as 128, and the dimension of edge embeddings n as 256. The weight of edge dictionary commitment loss  $\alpha$  is 0.25, and the weight of edge loss  $\beta$  is 0.1. The bins of event voxel grid B is 5. For training, we randomly apply color jitter, horizontal flipping, and gaussian blur to images and randomly resize with scales from 0.5 to 2.0 to images and events and crop the inputs to 256 × 256. For testing, we follow the setting of CMX [50] and CMNeXt [51], which upsamples the inputs to a width and height both divisible by 32 (will be ablated in table 3). More detailed information on training settings can be seen in appendix B.

Table 1: Compa	arisons on DERS-XS.	DERS-XR	DSEC-Semantic.	and DSEC-Xtrm.

Methods	Modality	DER		DERS		DSEC-S	emantic	DSEC	-Xtrm
Wichiods	wiodanty	mACC(%)↑	mIoU(%)↑	mACC(%)↑	mIoU(%)↑	mACC(%)↑	mIoU(%)↑	mACC(%)↑	mIoU(%)↑
SegFormer [47]	RGB	62.45	53.21	55.37	51.09	72.91	65.03	41.74	33.88
SegFormer (E) [47]	] Event	47.85	37.32	42.03	36.96	47.38	38.59	48.52	37.72
EvSegFormer [15]	Event	41.48	31.85	38.60	33.66	44.72	37.13	42.33	34.68
TokenFusion [41]	E-RGB	64.88	56.22	54.19	47.72	74.60	67.39	53.04	45.41
CMX [50]	E-RGB	71.86	63.12	64.51	59.22	76.18	68.10	51.29	43.95
CMNeXt [51]	E-RGB	73.30	64.55	66.57	60.95	77.50	69.03	52.12	45.16
EISNet [46] †	E-RGB	69.10	60.68	66.18	61.81	71.60	64.67	56.77	48.76
Ours	E-RGB	75.26	67.10	70.75	65.22	78.61	71.04	59.45	50.87

<sup>&</sup>lt;sup>†</sup> Reimplemented on DSEC-Semantic using the same dataloader for fair comparison with different training settings from [46], including each sequence events count of 50000 ([46] of 100000), input cropping size of 256 × 256 ([46] of 448 × 448), total batch size of 32 ([46] of 8), different random resize strategy, *etc*.

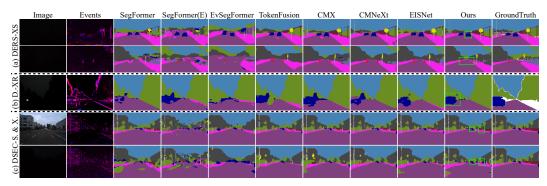


Figure 7: Qualitative comparison on DERS-XS, DERS-XR, DSEC-Semantic, and DSEC-Xtrm.

### 4.2 Comparisons with state-of-the-art

We compare our model with the current state-of-the-art, including RGB-only, event-only, and event-RGB-based methods. For RGB-only, we reimplement SegFormer [47] with MiT-B2 backbone settings, which is a powerful RGB-based semantic segmentation architecture. For event-only, we reimplement EvSegFormer [15] and modify the number of input channels of SegFormer to adapt to events input as an events-only comparison method. For event-RGB-based, we reimplement TokenFusion [41], CMX [50], CMNeXt [51], and EISNet [43]. For event-version SegFormer and EvSegFormer, we follow the setting of EvSegFormer and apply 6-channel image [1] as their event representation. For TokenFusion, CMX, CMNeXt, and EISNet, we use 3-bin voxel grid [61] as their event representation. We retrain all methods with the same training settings for fair comparison.

**Comparisons on true-labeled synthetic dataset.** As shown in table 1, our method surpasses all uni-modality and multi-modality methods and outperforms CMNeXt by a 2.55% mIoU on DERS-XS. As shown in fig. 7a, our method is more stable and robust to the edges of the segmentation results, especially for moving vehicles and pedestrians. This shows that our model effectively leverages edge information from events to compensate for the information loss of RGB under extreme conditions.

Comparisons on true-labeled real-world dataset fine-tuning. As shown in table 1, our methods outperforms EISNet by a 3.41% mIoU on DERS-XR fine-tuning experiments. Results also demonstrate that models trained on synthetic DERS-XS can be efficiently adapted to real-world data with minimal fine-tuning, further validating the effectiveness of DERS-XS. As shown in fig. 7b, our method successfully segments the vehicles, while other methods fail under real-world extreme scenes.

**Comparisons on non-extreme DSEC-Semantic.** As shown in table 1, our method outperforms CMNeXt by a 2.01% mIoU on DSEC-Semantic, demonstrating its effectiveness on a publicly available dataset under real-world, non-extreme conditions, despite the inherent limitations of pseudolabels. EISNet performs slightly worse, possibly due to its sensitivity to the input cropping strategy.

Comparisons under extreme conditions on degraded DSEC-Xtrm. As shown in table 1, our method ourperforms EISNet by a 2.11% mIoU on DSEC-Xtrm. Results also show that our method suffers less performance degradation with degraded inputs. As shown in fig. 7c, our method preserves more complete boundaries for vehicles and pedestrians, demonstrating its robustness and resilience.

#### 4.3 Ablation studies and analyses

Resilience study under severe spatial occlusion. This study emulates visual degradation due to spatial information loss under extreme conditions by applying local masking to the inputs. If masking is applied, we mask a  $100 \times 100$ area starting at coordinates (350, 200) for RGB and  $\langle 150, 150 \rangle$  for event. As shown in table 2, our method suffers less performance degradation under different settings. As shown in fig. 8, under E-RGB masking, CMX and CMNeXt fail in understanding the semantics of the mask area, while our method overcomes the problem by edge-aware optimization with uncertainty indicators. This demonstrates that our method is more resilient than other methods under severe spatial occlusion. Extended experiments with more results can be seen in appendix F and G.

Ablation study on ESC architecture. As shown in table 3, we ablate our ESC by removing modules under no mask setting and E-RGB mask setting on DERS-XS. Upsampling, as a form of data augmentation, is first ablated. Results show that each proposed module contributes positively to performance and resilience.

Ablation study on different K of edge dictionary. As shown in table 4, when K is too small, the items are insufficient for edge representation; when K is too large, the excess items are underutilized, leading to confusion in model learning. Both result in performance degradation; thus, we set K=128 to achieve an optimal trade-off.

Table 2: Study on spatial occlusion on DERS-XS.

Methods/mIoU(%)↑	Apply masking on								
Methods/fillioc(%)	None	RGB	Event	E-RGB					
TokenFusion [41]	56.22	48.44	55.79	48.00					
CMX [50]	63.12	54.13	62.70	53.73					
CMNeXt [51]	64.55	54.15	64.07	53.70					
EISNet [46]	60.68	55.33	59.87	54.47					
Ours	67.10	64.34	66.65	63.87					

Table 3: Ablation study on architecture on DERS-XS.

Arch./mIoU(%)↑	#Params(M)	w/o masking	E-RGB masking
ESC	56.875	67.10	63.87
<ul> <li>w/o Upsampling</li> </ul>	56.875	64.59	61.46
- w/o RC	56.612	64.29 (-0.31)	59.53 (-1.93)
- w/o UO	56.084	62.53 (-2.06)	58.43 (-3.03)
- w/o ELR& $L_{edge}$	38.411	61.35 (-3.24)	56.34 (-5.12)

Table 4: Ablation study on key usage on DERS-XS.

K	16	32	64	128	256	512
K-Usage <sup>†</sup>	16	32	64	92	99	97
gACC(%)↑	92.66	93.12	93.03	93.27	93.19	92.93
mACC(%)↑	74.59	74.77	74.90	75.26	73.91	73.65
mIoU(%)↑	65.87	66.84	66.64	67.10	66.54	66.11

<sup>† :</sup> K-Usage is the number of dictionary keys used.

Table 5: Ablation study on noise embeddings removal. Arch./mIoU(%) $\uparrow$  DSEC-XS DSEC-Semantic DSEC-Xtrm ESC w/o N<sub>K</sub>, N<sub>V</sub> 66.05 70.86 50.05 ESC w/ N<sub>K</sub>, N<sub>V</sub> 67.10 71.04 50.87

Table 6: Comp. on model complexity on DERS-XS.

Methods	Backbone	#Params(M)	FLOPs(G)	mIoU(%)↑
SegFormer [47]	MiT-B2	24.725	25.279	53.21
SegFormer (E) [47]	MiT-B2	24.734	25.433	37.32
EvSegFormer [15]	MiT-B2	24.740	25.422	31.85
TokenFusion [41]	MiT-B2	26.011	54.845	56.22
CMX [50]	$2 \times MiT-B2$	66.566	65.551	63.12
CMNeXt [51]	$2 \times MiT-B2$	58.687	62.805	64.55
EISNet [46]	MiT-B0 + B2	34.367	67.304	60.68
Ours	MiT-B1 + B2	56.875	95.086	67.10

Ablation study on the removal of noise embeddings. As shown in table 5, the removal of noise embeddings leads to a 1.05% and 0.81% mIoU drop on DERS-XS and DSEC-Xtrm, respectively, confirming their contribution to improved fitting stability. The performance drop on DSEC-Semantic is minimal (0.18%), which we attribute to its reliance on the RGB-based pseudo-labels. As the supervision signal has a bias towards RGB, the model naturally relies less on event modality. In such cases, the role of noise embeddings in facilitating cross-modality interaction becomes less significant.

**Model complexity.** Table 6 summarizes the complexity of compared models, where FLOPs are calculated with inputs of 512 × 512. Results show that our method has fewer parameters than CMX and CMNeXt, yet achieves better performance. The FLOPs of our method are relatively large, primarily due to multiple MLP heads for re-coding and resilient fusion in our framework.

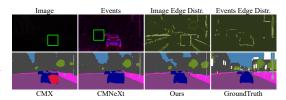


Figure 8: Qualitative study under spatial occlusion.

# 5 Conclusions and limitations

**Conclusions.** In this paper, we propose Edge-awareness Semantic Concordance, a multi-modality framework for event-RGB semantic segmentation. We demonstrate its capability and robustness for handling heterogeneous event and RGB. Results show that our framework outperforms existing event-RGB segmentation methods and possesses superior resilience in the case of modality imbalance and failure under extreme conditions. **Limitations.** Despite the promising results, only the fusion of event and RGB is considered so far. Exploring interactions with other visual modalities and designing modules tailored to their specific characteristics continues to be an open direction for future research.

# Acknowledgements

This work is partially supported by grants from the National Natural Science Foundation of China under contracts No. 62132002 and No. 62202010, the Beijing Nova Program (No.20250484786), and the Fundamental Research Funds for the Central Universities.

### References

- [1] Inigo Alonso and Ana C Murillo. Ev-segnet: Semantic segmentation for event-based cameras. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pages 0–0, 2019.
- [2] Bryce Bayer. Color imaging array. United States Patent, no. 3971065, 1976.
- [3] Yoshua Bengio, Nicholas Léonard, and Aaron Courville. Estimating or propagating gradients through stochastic neurons for conditional computation. *arXiv* preprint arXiv:1308.3432, 2013.
- [4] Jonathan Binas, Daniel Neil, Shih-Chii Liu, and Tobi Delbruck. Ddd17: End-to-end davis driving dataset. *arXiv preprint arXiv:1711.01458*, 2017.
- [5] Shristi Das Biswas, Adarsh Kosta, Chamika Liyanagedera, Marco Apolinario, and Kaushik Roy. Halsie: Hybrid approach to learning segmentation by simultaneously exploiting image and event modalities. In 2024 IEEE/CVF Winter Conference on Applications of Computer Vision (WACV), pages 5952–5962. IEEE, 2024.
- [6] Christian Brandli, Raphael Berner, Minhao Yang, Shih-Chii Liu, and Tobi Delbruck. A 240 × 180 130 db 3 μs latency global shutter spatiotemporal vision sensor. *IEEE Journal of Solid-State Circuits*, 49(10):2333–2341, 2014.
- [7] Tim Brooks, Ben Mildenhall, Tianfan Xue, Jiawen Chen, Dillon Sharlet, and Jonathan T Barron. Unprocessing images for learned raw denoising. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11036–11045, 2019.
- [8] Jiaqi Chen, Jiachen Lu, Xiatian Zhu, and Li Zhang. Generative semantic segmentation. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 7111–7120, 2023.
- [9] Alexey Dosovitskiy, German Ros, Felipe Codevilla, Antonio Lopez, and Vladlen Koltun. CARLA: An open urban driving simulator. In *Proceedings of the 1st Annual Conference on Robot Learning*, pages 1–16, 2017.
- [10] Guillermo Gallego, Tobi Delbrück, Garrick Orchard, Chiara Bartolozzi, Brian Taba, Andrea Censi, Stefan Leutenegger, Andrew J Davison, Jörg Conradt, Kostas Daniilidis, et al. Event-based vision: A survey. *IEEE transactions on pattern analysis and machine intelligence*, 44(1):154–180, 2020.
- [11] Mathias Gehrig, Willem Aarents, Daniel Gehrig, and Davide Scaramuzza. Dsec: A stereo event camera dataset for driving scenarios. *IEEE Robotics and Automation Letters*, 6(3):4947–4954, 2021.
- [12] Mengyue Geng, Lin Zhu, Lizhi Wang, Wei Zhang, Ruiqin Xiong, and Yonghong Tian. Event-based visible and infrared fusion via multi-task collaboration. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 26929–26939, 2024.
- [13] Yuhuang Hu, Shih-Chii Liu, and Tobi Delbruck. v2e: From video frames to realistic dvs events. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 1312–1321, 2021.
- [14] Ding Jia, Jianyuan Guo, Kai Han, Han Wu, Chao Zhang, Chang Xu, and Xinghao Chen. Geminifusion: Efficient pixel-wise multimodal fusion for vision transformer. In *Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024*. OpenReview.net, 2024.
- [15] Zexi Jia, Kaichao You, Weihua He, Yang Tian, Yongxiang Feng, Yaoyuan Wang, Xu Jia, Yihang Lou, Jingyi Zhang, Guoqi Li, et al. Event-based semantic segmentation with posterior attention. *IEEE Transactions on Image Processing*, 32:1829–1842, 2023.
- [16] Yu Jiang, Yuehang Wang, Siqi Li, Yongji Zhang, Minghao Zhao, and Yue Gao. Event-based lowillumination image enhancement. IEEE Transactions on Multimedia, 2023.

- [17] Zhe Jiang, Yu Zhang, Dongqing Zou, Jimmy Ren, Jiancheng Lv, and Yebin Liu. Learning event-based motion deblurring. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3320–3329, 2020.
- [18] Linglin Jing, Yiming Ding, Yunpeng Gao, Zhigang Wang, Xu Yan, Dong Wang, Gerald Schaefer, Hui Fang, Bin Zhao, and Xuelong Li. Hpl-ess: Hybrid pseudo-labeling for unsupervised event-based semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 23128–23137, 2024.
- [19] Donggeun Kim and Taesup Kim. Missing modality prediction for unpaired multimodal learning via joint embedding of unimodal models. In *European Conference on Computer Vision*, pages 171–187. Springer, 2024.
- [20] Jian Lang, Zhangtao Cheng, Ting Zhong, and Fan Zhou. Retrieval-augmented dynamic prompt tuning for incomplete multimodal learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 18035–18043, 2025.
- [21] Hebei Li, Jin Wang, Jiahui Yuan, Yue Li, Wenming Weng, Yansong Peng, Yueyi Zhang, Zhiwei Xiong, and Xiaoyan Sun. Event-assisted low-light video object segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3250–3259, 2024.
- [22] Xiangtai Li, Xia Li, Li Zhang, Guangliang Cheng, Jianping Shi, Zhouchen Lin, Shaohua Tan, and Yunhai Tong. Improving semantic segmentation via decoupled body and edge supervision. In Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XVII 16, pages 435–452. Springer, 2020.
- [23] Jinxiu Liang, Yixin Yang, Boyu Li, Peiqi Duan, Yong Xu, and Boxin Shi. Coherent event guided low-light video enhancement. In Proceedings of the IEEE/CVF International Conference on Computer Vision, pages 10615–10625, 2023.
- [25] Alex Liu, SouYoung Jin, Cheng-I Lai, Andrew Rouditchenko, Aude Oliva, and James Glass. Cross-modal discrete representation learning. In Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 3013–3035, 2022.
- [26] Lin Liu, Junfeng An, Jianzhuang Liu, Shanxin Yuan, Xiangyu Chen, Wengang Zhou, Houqiang Li, Yan Feng Wang, and Qi Tian. Low-light video enhancement with synthetic event guidance. In *Proceedings* of the AAAI Conference on Artificial Intelligence, pages 1692–1700, 2023.
- [27] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In 7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019. OpenReview.net, 2019.
- [28] Jia Ning, Chen Li, Zheng Zhang, Chunyu Wang, Zigang Geng, Qi Dai, Kun He, and Han Hu. All in tokens: Unifying output space of visual tasks via soft token. In *Proceedings of the IEEE/CVF International* Conference on Computer Vision, pages 19900–19910, 2023.
- [29] Liyuan Pan, Cedric Scheerlinck, Xin Yu, Richard Hartley, Miaomiao Liu, and Yuchao Dai. Bringing a blurry frame alive at high frame-rate with an event camera. In *Proceedings of the IEEE/CVF Conference* on Computer Vision and Pattern Recognition, pages 6820–6829, 2019.
- [30] Yunshan Qi, Lin Zhu, Yu Zhang, and Jia Li. E2nerf: Event enhanced neural radiance fields from blurry images. In Proceedings of the IEEE/CVF International Conference on Computer Vision, pages 13254– 13264, 2023.
- [31] Henri Rebecq, René Ranftl, Vladlen Koltun, and Davide Scaramuzza. High speed and high dynamic range video with an event camera. *IEEE transactions on pattern analysis and machine intelligence*, 43(6):1964–1980, 2019.
- [32] Wei Shang, Dongwei Ren, Dongqing Zou, Jimmy S Ren, Ping Luo, and Wangmeng Zuo. Bringing events into video deblurring with non-consecutively blurry frames. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4531–4540, 2021.
- [33] Peilun Shi, Jiachuan Peng, Jianing Qiu, Xinwei Ju, Frank Po Wen Lo, and Benny Lo. Even: An event-based framework for monocular depth estimation at adverse night conditions. In 2023 IEEE International Conference on Robotics and Biomimetics (ROBIO), pages 1–7. IEEE, 2023.

- [34] Leslie N Smith. Cyclical learning rates for training neural networks. In 2017 IEEE winter conference on applications of computer vision (WACV), pages 464–472. IEEE, 2017.
- [35] Timo Stoffregen, Cedric Scheerlinck, Davide Scaramuzza, Tom Drummond, Nick Barnes, Lindsay Kleeman, and Robert Mahony. Reducing the sim-to-real gap for event cameras. In Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXVII 16, pages 534–549. Springer, 2020.
- [36] Zhaoning Sun, Nico Messikommer, Daniel Gehrig, and Davide Scaramuzza. Ess: Learning event-based semantic segmentation from still images. In *European Conference on Computer Vision*, pages 341–357. Springer, 2022.
- [37] Andrew Tao, Karan Sapra, and Bryan Catanzaro. Hierarchical multi-scale attention for semantic segmentation. arXiv preprint arXiv:2005.10821, 2020.
- [38] Gemma Taverni, Diederik Paul Moeys, Chenghan Li, Celso Cavaco, Vasyl Motsnyi, David San Segundo Bello, and Tobi Delbruck. Front and back illuminated dynamic and active pixel vision sensors comparison. *IEEE Transactions on Circuits and Systems II: Express Briefs*, 65(5):677–681, 2018.
- [39] Aaron Van Den Oord, Oriol Vinyals, et al. Neural discrete representation learning. *Advances in neural information processing systems*, 30, 2017.
- [40] Lin Wang, Yujeong Chae, Sung-Hoon Yoon, Tae-Kyun Kim, and Kuk-Jin Yoon. Evdistill: Asynchronous events to end-task learning via bidirectional reconstruction-guided cross-modal knowledge distillation. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 608–619, 2021.
- [41] Yikai Wang, Xinghao Chen, Lele Cao, Wenbing Huang, Fuchun Sun, and Yunhe Wang. Multimodal token fusion for vision transformers. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 12186–12195, 2022.
- [42] Wentao Wu, Xiao Wang, Chenglong Li, Bo Jiang, Jin Tang, Bin Luo, and Qi Liu. Cm3ae: A unified rgb frame and event-voxel/-frame pre-training framework. *arXiv preprint arXiv:2504.12576*, 2025.
- [43] Ruihao Xia, Chaoqiang Zhao, Meng Zheng, Ziyan Wu, Qiyu Sun, and Yang Tang. CMDA: cross-modality domain adaptation for nighttime semantic segmentation. In *IEEE/CVF International Conference on Computer Vision, ICCV 2023, Paris, France, October 1-6, 2023*, pages 21515–21524. IEEE, 2023.
- [44] Yan Xia, Hai Huang, Jieming Zhu, and Zhou Zhao. Achieving cross modal generalization with multimodal unified representation. *Advances in Neural Information Processing Systems*, 36, 2024.
- [45] Xiaoyang Xiao, Yuqian Zhao, Fan Zhang, Biao Luo, Lingli Yu, Baifan Chen, and Chunhua Yang. Baseg: Boundary aware semantic segmentation for autonomous driving. *Neural Networks*, 157:460–470, 2023.
- [46] Bochen Xie, Yongjian Deng, Zhanpeng Shao, and Youfu Li. Eisnet: A multi-modal fusion network for semantic segmentation with events and images. *IEEE Transactions on Multimedia*, 2024.
- [47] Enze Xie, Wenhai Wang, Zhiding Yu, Anima Anandkumar, Jose M Alvarez, and Ping Luo. Segformer: Simple and efficient design for semantic segmentation with transformers. *Advances in neural information processing systems*, 34:12077–12090, 2021.
- [48] Yan Yang, Liyuan Pan, and Liu Liu. Event camera data pre-training. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 10699–10709, 2023.
- [49] Bowen Yao, Yongjian Deng, Yuhan Liu, Hao Chen, Youfu Li, and Zhen Yang. Sam-event-adapter: Adapting segment anything model for event-rgb semantic segmentation. In 2024 IEEE International Conference on Robotics and Automation (ICRA), pages 9093–9100. IEEE, 2024.
- [50] Jiaming Zhang, Huayao Liu, Kailun Yang, Xinxin Hu, Ruiping Liu, and Rainer Stiefelhagen. Cmx: Cross-modal fusion for rgb-x semantic segmentation with transformers. *IEEE Transactions on intelligent transportation systems*, 2023.
- [51] Jiaming Zhang, Ruiping Liu, Hao Shi, Kailun Yang, Simon Reiß, Kunyu Peng, Haodong Fu, Kaiwei Wang, and Rainer Stiefelhagen. Delivering arbitrary-modal semantic segmentation. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 1136–1147, 2023.
- [52] Jiaming Zhang, Kailun Yang, and Rainer Stiefelhagen. Issafe: Improving semantic segmentation in accidents by fusing event-based data. In 2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS), pages 1132–1139. IEEE, 2021.

- [53] Song Zhang, Yu Zhang, Zhe Jiang, Dongqing Zou, Jimmy Ren, and Bin Zhou. Learning to see in the dark with events. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XVIII 16*, pages 666–682. Springer, 2020.
- [54] Fengan Zhao, Qianang Zhou, and Junlin Xiong. Edge-guided fusion and motion augmentation for event-image stereo. In European Conference on Computer Vision, pages 190–205. Springer, 2024.
- [55] Yucheng Zhao, Gengyu Lyu, Ke Li, Zihao Wang, Hao Chen, Zhen Yang, and Yongjian Deng. Eseg: Event-based segmentation boosted by explicit edge-semantic guidance. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 10510–10518, 2025.
- [56] Mingmin Zhen, Jinglu Wang, Lei Zhou, Shiwei Li, Tianwei Shen, Jiaxiang Shang, Tian Fang, and Long Quan. Joint semantic segmentation and boundary detection using iterative pyramid contexts. In *Proceedings* of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 13666–13675, 2020.
- [57] Sipeng Zheng, Bohan Zhou, Yicheng Feng, Ye Wang, and Zongqing Lu. Unicode: Learning a unified codebook for multimodal large language models. *arXiv preprint arXiv:2403.09072*, 2024.
- [58] Chu Zhou, Minggui Teng, Jin Han, Jinxiu Liang, Chao Xu, Gang Cao, and Boxin Shi. Deblurring low-light images with events. *International Journal of Computer Vision*, 131(5):1284–1298, 2023.
- [59] Chu Zhou, Minggui Teng, Jin Han, Chao Xu, and Boxin Shi. Delieve-net: Deblurring low-light images with light streaks and local events. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1155–1164, 2021.
- [60] Hanyu Zhou, Yi Chang, Haoyue Liu, Wending Yan, Yuxing Duan, Zhiwei Shi, and Luxin Yan. Exploring the common appearance-boundary adaptation for nighttime optical flow. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net, 2024.
- [61] Alex Zihao Zhu, Liangzhe Yuan, Kenneth Chaney, and Kostas Daniilidis. Unsupervised event-based optical flow using motion compensation. In *Proceedings of the European Conference on Computer Vision* (ECCV) Workshops, pages 711–714, 2018.

# **Appendix**

# A Details of datasets

Three extreme-condition datasets are constructed and used for testing in our work, namely DERS-XS, DERS-XR, and DSEC-Xtrm. Among them, DERS-XS and DSEC-Xtrm are synthetic datasets, and DERS-XR is a real-world dataset. In addition, a common dataset of normal conditions, namely DSEC-Semantic is used for testing. By leveraging edge-awareness, our method can effectively obtain the common features of heterogeneous event and RGB under unified semantic space and jointly optimize them. Results show that our method outperforms existing event-RGB segmentation methods and possesses superior resilience in the case of modality imbalance and failure under extreme conditions.

Datasets #Train #Validation #Test Real/Syn. True/Pseudo-lbl. Fine/Coarse-lbl. Avg. Pixel Val. Avg. #Events DERS-XS 2016 1080 70490.88 144 Synthetic Fine-label 6.16 True-label DERS-XR 120 N/A 120 Real-world True-label Coarse-label 20.58 12784 76 DSEC-Semantic 8082 N/A 2809 Real-world Pseudo-label Fine-label 75 13 608162.26 2809 DSEC-Xtrm 8082 N/A Synthetic Pseudo-label Fine-label 4.75 689098.25

Table 7: Comparison between datasets.

## A.1 DERS-XS

Our synthetic dataset DERS-XS is constructed based on the CARLA simulator [9] and v2e simulator [13], containing 270 frame sequences. We first obtain canonical RGB frames with labels from CARLA simulator. For CARLA simulation process, we first load six pre-made maps, namely *Town01*, *Town02*, *Town03*, *Town04*, *Town05*, *and Town10*. Then we apply fifteen different types of weather conditions on each map. For each weather and map combination, we record two sequences of 1200 frames at a frame rate of 20 fps. The spatial size of each frame is  $640 \times 360$ .

The CARLA simulated segmentation labels contain 23 categories, and the specific category names are unlabeled, building, fence, other, pedestrian, pole, roadline, road, sidewalk, vegetation, vehicles, wall, traffic sign, sky, ground, bridge, rail track, guard rail, traffic light, static, dynamic, water, and terrain. For experiments, we merge and transform the above 23 categories into 11 categories, in order to match the categories setting of DSEC-Semantic [11, 36]. The 11 categories are background, building, fence, person, pole, road, sidewalk, vegetation, car, wall, and traffic sign. Categories that do not exist in DSEC-Semantic are set to 255 and ignored during the training and testing process.

We then obtain noisy events from v2e simulator based on the CARLA simulated frames. For v2e simulation process, the positive threshold and the negative threshold are both set as 0.2. The 1-std deviation threshold variation is set as 0.05. The cutoff frequency is set as 30, and the leak event rate per pixel is set as 0.1. The shot noise rate is set as 5.0. The refractory period is set as 0.0005.

We implement the Image Signal Processing (ISP) pipeline and its inverse process, referred to as the ISP unprocessing technique in [7], which includes digital gain, white balance, demosaicing, color correction, gamma compression, and tone mapping. By utilizing the inversion of ISP, we first convert CARLA-simulated canonical RGB frames into Bayer-pattern BGGR RAW images [2]. As shown in fig. 9, we attenuate the optical signals and add shot noise on RAW images, and then process RAW images by ISP to obtain low-light images.

We sample the 1200-frame sequences at the intervals of 100 frames, taking 12 frames from each sequence while discarding the rest. We divide the 270 sequences into three parts, of which the training set has 168 sequences, the validation set has 12 sequences, and the test set has 90 sequences. For training process of DERS-XS, we use the training set for training, and the validation set for saving the model with the best validation mIoU. The test set is used only during the testing process.

### A.2 DERS-XR

Our real-world dataset DERS-XR is captured by a DAVIS-346 [38], with the spatial size of 346 × 260. We use the camera to capture the events with APS frames of 20 fps simultaneously under extreme lighting conditions. We capture a total of 27 frame sequences, sampling them at intervals

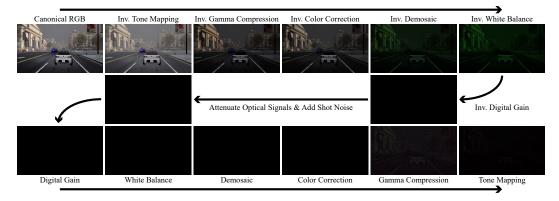


Figure 9: **Obtaining low-light RGB images using ISP with its inversion process.** We implement the ISP with its inversion process, including digital gain, white balance, demosaic, color correction, gamma compression, and tone mapping. Optical Signals are attenuated, and shot noise is added on RAW domain images. The process is used to simulate low-light RGB images in DERS-XS and DSEC-Xtrm.

of 100 frames, and finally obtain 240 frames of images with very different contents. We manually annotate semantic labels for the sampled 240 frames, and the annotated categories are the same as the transformed 11 categories of DSEC-XS. We randomly select 120 frames for fine-tuning and the remaining 120 frames for testing. For experiments conducted on DERS-XR, we save the last epoch fine-tuning model for testing.

#### A.3 DSEC-Xtrm

DSEC-Xtrm is simulated and converted from the real-world dataset DSEC-Semantic. A sample of DSEC-Semantic and DSEC-Xtrm is shown in fig. 10. DSEC-Semantic includes 11 sequences with pseudo-labels of 19 categories and 11 categories. We generate pure noise events by modifying the source codes of v2e simulator and overlaying them to the event sequence of DSEC-Semantic to obtain noisy events in DSEC-Xtrm. The shot noise rate is set as 10.0 Hz. To obtain low-light images, we apply the same ISP and ISP-inversion process, unprocessing RGB images to RAW domain, attenuating the optical signals and adding shot noise on the unprocessed RAW images, and then process the RAW images to RGB low-light images by ISP. We sample the 11 sequences at the intervals of 2 frames and discard the first 6 frames of each sequence, fol-

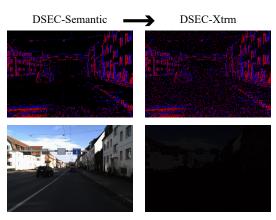


Figure 10: DSEC-Semantic and DSEC-Xtrm.

lowing the same setting with [36]. The simulated low-light images and noisy events combined with the original 11-categories labels of DSEC-Semantic together constitute the final DSEC-Xtrm dataset.

## A.4 Discussion on datasets

This paper uses the DSEC-Semantic and constructs three datasets, each of which has its own properties. Table 7 compares the properties of different datasets, including the number of image-events-pair in the training set, validationset and test set, whether the data is real-world or synthetic, whether the label is true or pseudo, whether the label is fine or coarse, and the average pixel value of images and average number of events of each dataset. Based on the different properties of different datasets, we can use them in different experimental settings to test the model comprehensively.

The DERS-XS only contains synthetic data, but DERS-XS has the largest amount of data with the true fine-grained label, thus DERS-XS can be a standard benchmark for comparative

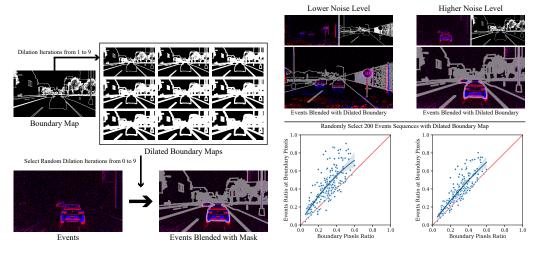


Figure 11: Details of boundary map dilation.

Figure 12: Event-edge statistics under two noise levels.

**experiments**. The DERS-XR is a real-world dataset with the most realistic data distribution, however, it is difficult to annotate it accurately, and for this reason, we only coarsely annotate a small amount of data of DERS-XR. Since it is difficult to annotate DERS-XR, only a small amount of data is annotated, thus we only use DERS-XR for fine-tuning and testing.

The DSEC-Semantic is a real-world dataset under normal conditions, however, the labels are pseudo labels based on RGB images only. Although it is used for testing in many works, **using DSEC-Semantic as a benchmark for multi-modality semantic segmentation is defective**. Thus, we simulate the extreme version DSEC-Xtrm from DSEC-Semantic. For DSEC-Xtrm, the pseudo labels from DSEC-Semantic are no longer a defect for being a benchmark for multi-modality semantic segmentation. **Therefore, when testing on the DSEC-Semantic and DSEC-Xtrm, we can compare the performance degradation on the two datasets for different methods.** The performance degradation on DSEC-Semantic and DSEC-Xtrm can illustrate the robustness and resilience of different methods.

# **B** Details of training settings

For all datasets, we use AdamW [27] optimizer with a weight decay of 0.01, and the learning rate of decoder is 10 times the basic learning rate. For DERS-XS, we train our model on two NVIDIA RTX 3090 GPUs for 300 epochs, and the batch size is 16 on each GPU. The basic learning rate (LR) is 6  $\times$  10<sup>-5</sup>, which is scheduled by a CyclicLR [34] scheduler with a maximum learning rate 1.6  $\times$  LR and a triangular cycle of 10 epochs. For DERS-XR, we fine-tune our model for 50 epochs on a single GPU with a batch size of 2, based on the best model trained on DERS-XS. The basic learning rate is 6  $\times$  10<sup>-5</sup>, which is scheduled by a WarmupPolyLR scheduler with power of 0.9 and warmup epochs of 10. For DSEC-Semantic and DSEC-Xtrm, we follow settings from DERS-XS with epochs of 60. We input events with a 50 ms interval for DERS-XS and DERS-XR, and a count of 50000 events for DSEC-Semantic and DSEC-Xtrm.

# C Details of event-edge statistics

For the statistical process, we randomly select 200 event sequences of 50 ms accompanied by their boundary maps from DERS-XS as sample. We aim to count the ratio of edge pixels to all pixels of the whole plane, and the ratio of events falling on edge pixels to all events of the whole sequence. As shown in fig. 11, in order to make the ratio of edge pixels distributed in a larger range, we first dilate the boundary map with a  $3 \times 3$  kernel with a random number of iterations in the range of 10. We draw a scatter plot to show the correlation between the two ratios.

As the boundary map dilates, the ratio of edge pixels increases, and the ratio of events falling on edge pixels also increases synchronously. Considering the extreme cases, when the ratio of edge pixels is

0, there is no events falling on edge pixels, then the ratio of events is also 0; when the boundary map dilates to the whole image, the ratio of edge pixels is 1, then all events fall on edge pixels, the ratio of events is 1. However, for the cases where the ratio of edge pixels is between 0 and 1, for most samples, the ratio of events falling on edge pixels is greater than the ratio of edge pixels. This proves the fact that events tend to cluster at the areas of semantic edge. This phenomenon exhibits a strong correlation between events and semantic edge.

We also count the two ratios under two noise levels of events. We use v2e simulator to simulate a version of event sequences that are less noisy than the event sequences of DERS-XS. We simultaneously count the event ratios corresponding to the two noise versions of the event sequences, and draw scatter plots for the lower noise level event sequences and the higher noise level event sequences (*i.e.* DERS-XS). We compare the scatter plots of the two noise levels. As shown in fig. 12, as the noise level increases, the ratio of edge pixels to all pixels and the ratio of events falling on edge pixels to all events tend to be equal. The correlation curve for the higher noise level case is still a concave curve, thus the correlation is still maintained even under the influence of high noise. This ensures the resilience of our method in the case of modality imbalance or failure under extreme conditions.

# D Details of edge dictionary training process

The training of our edge dictionary is a separate stage from the training of the segmentation model. Through this separate training stage, we obtain the pre-trained weights of edge dictionary, which represents the information of the semantic edge. We also obtain the pre-trained weights of the tokenizer, which can be used to embed the semantic edge ground truth into the discrete latent space defined by edge dictionary. The detokenizer is only used in the edge dictionary training stage, and is deprecated in the segmentation model training stage.

The tokenizer  $f_T$  consists of two convolutional layers followed by ReLU, two residual blocks, and a final convolutional layer. The first two convolutional layers downsampled the inputs each with a kernel of  $4 \times 4$  size, a stride of  $\langle 2, 2 \rangle$ , and a padding of  $\langle 1, 1 \rangle$ . The two residual blocks keep the spatial scale unchanged, each consists three convolutional layers followed by ReLU with a kernel of  $3 \times 3$  size, a stride of  $\langle 1, 1 \rangle$ , and a padding of  $\langle 1, 1 \rangle$ . The final convolutional layers adjust the number of channels to n by a convolutional layer with a kernel of  $1 \times 1$  size. For semantic edge  $\mathbf{B} \in \{0,1\}^{H \times W}$ , the final produced edge embeddings  $\mathbf{\Gamma} = f_T(\mathbf{B}) \in \mathbb{R}^{H' \times W' \times n}$  have the downsampled spatial size  $H' \times W'$ , where  $H' = \lfloor \frac{H}{4} \rfloor$ ,  $W' = \lfloor \frac{W}{4} \rfloor$ .

The detokenizer  $f_{T'}$  is constructed by changing the downsampled convolutional layers to the transposed convolutional layer for upsampling and inverting all the layers in tokenizer. Thus, the detokenizer consists of a convolutional layer followed by ReLU, two residual blocks and two transposed convolutional layers followed by ReLU, and a final convolutional layer to predict the semantic edge. The kernel size, stride, and padding settings of convolutional layers and transposed convolutional layers are the same as the corresponding layers in tokenizer. The detokenizer takes the quantised embeddings  $\Gamma' \in \mathbb{R}^{H' \times W' \times n}$  as input, predicts the reconstructed semantic edge  $\mathbf{B}' = f_{T'}(\Gamma')$ .

We adopt the training objective with reconstruction loss, embedding loss, and commitment loss of VQ-VAE [39] as  $L_{dict} = \|\mathbf{B} - \mathbf{B}'\|_2^2 + \|v(\hat{\mathbf{K}}) - \mathrm{sg}(\Gamma)\|_2^2 + \alpha\|\mathrm{sg}(v(\hat{\mathbf{K}})) - \Gamma\|_2^2$ , where sg means stop gradient, and  $\alpha$  is a constant of commitment loss weight, which is 0.25 in our work. To make the reconstruction loss propagate back to the tokenizer, a gradient straight-through technique [3] is adopted, which directly assigns the gradient from  $\Gamma'$  to  $\Gamma$ . We train the edge dictionary with tokenizer and detokenizer for 3000 epochs, based on the data of DERS-XS and DSEC-Semantic. No additional information is introduced, and the pre-trained weights of tokenizer are introduced into the segmentation training stage only for the supervision in latent space, and not utilized for testing stage.

### E Additional ablation studies

## E.1 Edge Dictionary Domain Transferability

Theoretically, the discrete edge dictionary learned by VQ-VAE is expected to have good transferability across datasets. As an intermediate representation, semantic edge exhibits relatively simple and consistent structures, and the latent distributions of semantic edge derived from segmentation labels

Table 8: Ablation study on edge dictionary domain transferability.

Settings	gACC(%)↑	mACC(%)↑	mIoU(%)↑
ESC on DERS-XS (w/ edge dictionary of DSEC)	93.23	74.96	66.44
ESC on DERS-XS	93.27	75.26	67.10
ESC on DSEC-Semantic (w/ edge dictionary of DERS-XS)	94.91	78.04	70.93
ESC on DSEC-Semantic	94.85	78.61	71.04
ESC on DSEC-Xtrm (w/ edge dictionary of DERS-XS)	88.57	58.00	50.65
ESC on DSEC-Xtrm	88.18	59.45	50.87

Table 9: Ablation study on different event sampling strategies.

Settings	gACC(%)↑	mACC(%)↑	mIoU(%)↑
ESC on DSEC-Semantic (50 ms)	94.76	78.44	70.83
ESC on DSEC-Semantic (100,000 events)	94.87	78.00	70.84
ESC on DSEC-Semantic (50,000 events)	94.85	78.61	71.04

tend to vary only slightly across different datasets. Therefore, we expect the performance degradation under dictionary transfer settings to be minimal.

We conduct cross-domain transferability evaluations by (i) using an edge dictionary pretrained on DSEC to evaluate on DERS-XS, and (ii) using a dictionary pretrained on DERS-XS to evaluate on DSEC-Semantic and DSEC-Xtrm. As shown in table 8, the performance drops are small in all cases. Specifically, the mIoU on DERS-XS drops by 0.66%, on DSEC-Semantic by 0.10%, and on DSEC-Xtrm by 0.21%, compared to the original non-exchanged dictionary settings. These results suggest that the learned edge dictionary generalizes well across domains, and our method remains robust under moderate domain shifts.

# **E.2** Event Sampling Strategy

The DSEC-Semantic event input is built by sampling a fixed number of events per voxel grid rather than a fixed time window, which follows the setting of ESS [36]. In ESS, the event input is built with 100,000 events per voxel grid. We found that a 50-ms fixed time window or 100,000 fixed number of events is relatively large, which reduces the performance of data preprocessing, and may lead to insufficient edge characteristic representation. After the above trade-offs, we decided to use 50,000 events per voxel grid for DSEC-Semantic as our event sampling strategy in our work.

To further demonstrate the impact of different event sampling strategies, we conduct experiments on DSEC-Semantic with a fixed time window of 50 ms, a fixed number of events of 100,000 per voxel grid, compared with 50,000 events per voxel grid in the main paper. As shown in table 9, different event sampling strategies present comparable results, with mIoU slightly lower (0.20% and 0.19% respectively) than the fixed number of events of 50,000 in the main paper.

# E.3 Deployment Efficiency and FLOPs Ablation

We conduct additional comparative experiments on DERS-Xtrm using smaller backbones (2× MiT-B0), reducing the FLOPs of ESC to be even lower than CMNeXt. In addition, we measure the end-to-end inference latency of CMNeXt and our ESC (both standard ESC and reduced variant) on a single NVIDIA GeForce RTX 3090 GPU with a batch size of 1. All latency measurements are performed with a fixed input size of 512× 512, with each measurement calculating the average execution time over 100 inferences, and we repeat 3 times for stability.

As shown in table 10, the reduced ESC variant still outperforms CMNeXt, achieving 49.06% mIoU vs. 45.16%, despite lower FLOPs (60.658G vs. 62.805G) and significantly fewer parameters (14.184M vs. 58.687M). This suggests that the performance gains stem from architectural design rather than merely an increased computational cost. Furthermore, as shown in table 10, the reduced ESC variant has an average inference latency of 29.37 ms, which is shorter than that of CMNeXt (29.79 ms), demonstrating its potential for more efficient deployment.

Above results indicate that even with lightweight backbones, our model maintains strong performance with higher inference speed, highlighting the effectiveness of our design beyond raw FLOPs. This reflects a favorable trade-off between efficiency and performance, which is essential for practical deployment in real-world systems.

Table 10: Ablation study on lighter backbones on DSEC-Xtrm.

Settings	gACC(%)↑	mACC(%)↑	mIoU(%)↑	mIoU(%)↑ #Params(M)		Latency (ms)				
Settings	gACC(%)	$ \mathcal{E}(\mathcal{E}(\mathcal{E}))  =  \mathcal{E}(\mathcal{E}(\mathcal{E}))  =  \mathcal{E}(\mathcal{E}(\mathcal{E})) $	FLOPs(G)	#1	#2	#3	Avg.			
CMNeXt [51]	87.04	52.12	45.16	58.687	62.805	29.75	29.78	29.83	29.79	
ESC (Reduced)	88.03	56.31	49.06	14.184	60.658	29.30	29.38	29.43	29.37	
ESC (Standard)	88.18	59.45	50.87	56.875	95.086	34.56	34.46	34.78	34.60	

Table 11: Extended experiments under severe spatial occlusion.

Methods / (%)↑	$50 \times 50$			-	100 × 100			150 × 150			$200 \times 200$			250 × 250	
` / 1					mACC										
TokenFusion [41]	88.03	62.54	52.50	84.92	59.90	48.00	79.63	56.34	43.85	75.25	52.08	39.76	73.27	49.68	37.39
CMX [50]	90.53	70.73	57.76	86.43	67.40	53.73	82.24	64.12	49.59	80.21	59.21	45.96	79.33	56.81	43.79
CMNeXt [51]	91.22	71.91	60.65	87.53	68.71	53.70	81.62	64.30	48.01	77.08	59.51	43.71	74.37	56.42	41.12
EISNet [46]	90.43	68.12	57.89	88.20	66.19	54.47	83.83	63.31	49.30	80.99	60.71	45.92	78.83	58.78	43.51
Ours	92.46	74.08	64.53	92.20	72.91	63.87	91.05	70.40	61.32	89.01	66.90	58.02	87.26	64.06	54.95

# F Extended experiments under severe spatial occlusion

We extend our experiments under severe spatial occlusion. In the main text, we apply masking on RGB and event with mask areas size of  $100 \times 100$ . We further conduct experiments with masking areas of different sizes. Masking areas of size  $50 \times 50$ ,  $150 \times 150$ ,  $200 \times 200$ , and  $250 \times 250$  are applied at coordinate  $\langle 350, 200 \rangle$  for RGB and  $\langle 150, 150 \rangle$  for event respectively. Both event and RGB are applied masking. The excess part is ignored if the masking area exceeds the spatial area.

Table 11 demonstrates the comparison results of extended experiments under severe spatial occlusion. As the size of masking area increases, the performance of all methods degrades, and our method consistently outperforms other multi-modality methods on different masking settings. CMX has a lower performance than CMNeXt when the masking areas are small, but it outperforms CMNeXt when the masking areas become larger. Figure 13 demonstrates the qualitative comparison results of extended experiments on masking settings on DERS-XS. Under different masking settings, although all methods are affected by the modality imbalance and information loss caused by masking, our method obtains more reliable information based on uncertainty edge-aware joint optimization and edge consolidation, thus avoiding the misleading information of masking areas as much as possible.

# **G** More qualitative comparison results

This section is an extension of the qualitative results of the experiments in the main text. Figures are placed at the end of the appendix. Figure 14 demonstrates more qualitative comparison results on DERS-XS. Figure 15 demonstrates more qualitative comparison results on DERS-XR. Figure 16 demonstrates more qualitative comparison results on DSEC-Semantic and DSEC-Xtrm. Figure 17 demonstrates more qualitative comparison results under E-RGB mask setting on DERS-XS.

As shown in fig. 14, compared with other methods, our method can segment vehicles and pedestrians with complex contours more stably and robustly under extreme conditions. Especially for pedestrians, the contours of pedestrians segmented by other methods are not sharp enough, and sometimes even fail to detect the existence of pedestrians. our method can effectively locate pedestrians and segment the contours of pedestrians accurately.

As shown in fig. 15, our method performs well on real-world extreme scenes when fine-tuned with a small amount of real-world data. This also confirms that our simulated dataset DERS-XS can effectively provide prior knowledge when the amount of real-world data is small.

As shown in fig. 16, our method can handle Event-RGB semantic segmentation under normal conditions well, and when the modality information is lost under extreme conditions, our method is still able to identify the contours of ambiguous vehicles and pedestrians and maintains the segmentation performance better than other methods.

As shown in fig. 17, our method can still achieve relatively accurate segmentation results even when the inputs are partially masked. The image edge distribution and events edge distribution indicate their different awareness of edges. The two modalities complement each other and are dynamically fused based on their confidences and uncertainties. The results show that our method is more robust and resilient than other methods in the cases of modality imbalance and failure.

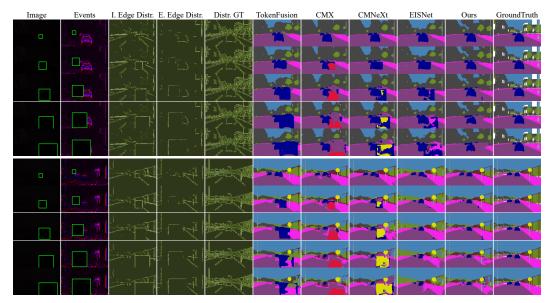


Figure 13: Qualitative results of extended experiments under severe spatial occlusion on DERS-XS.

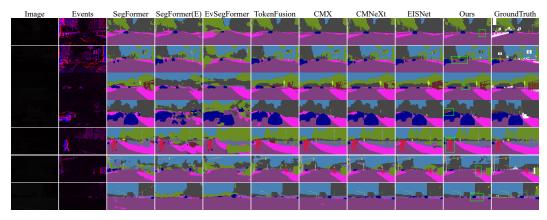


Figure 14: More qualitative comparison results on DERS-XS.

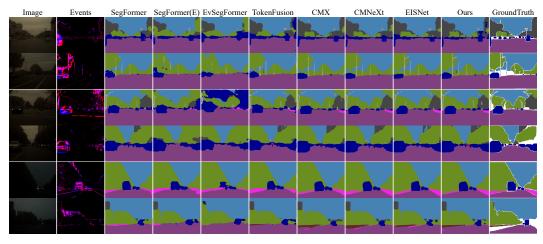


Figure 15: More qualitative comparison results on DERS-XR.

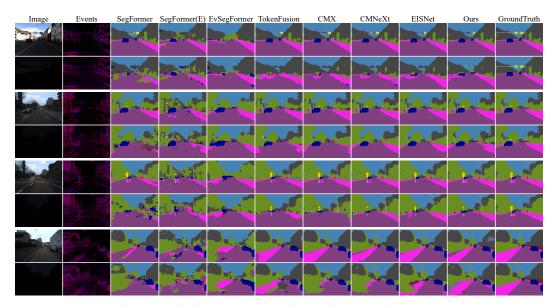


Figure 16: More qualitative comparison results on DSEC-Semantic and DSEC-Xtrm.

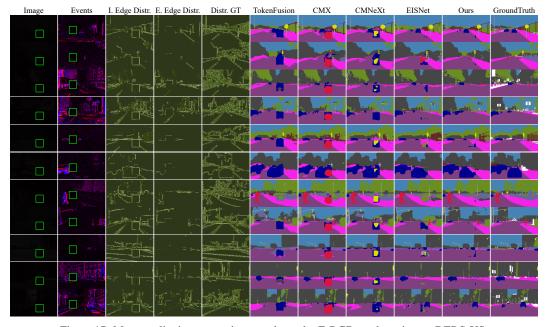


Figure 17: More qualitative comparison results under E-RGB mask setting on DERS-XS.

# **NeurIPS Paper Checklist**

### 1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: We are convinced that the main claims made in the abstract and introduction do accurately reflect the paper's contributions and scope.

#### Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals
  are not attained by the paper.

## 2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: Despite the limited length of the paper, we still briefly discussed one possible limitation of the work in section 5.

#### Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

### 3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [NA]

Justification: The results of the paper are all experimental, not theoretical.

### Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and crossreferenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

# 4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: Technical details of the core method and dataset construction method are provided in section 3, and implementation details are listed in section 4.1 and appendix B.

#### Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
  - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
- (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
- (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
- (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

# 5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: The code and datasets are released at http://github.com/iCVTEAM/ESC. Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

### 6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: All the training and test details are listed in section 4.1 and appendix B.

# Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental
  material.

# 7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [No]

Justification: Error bars are not reported due to the high computational cost.

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)

- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

# 8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: The type and amount of GPU are reported in section 4.1.

#### Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

### 9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?

Answer: [Yes]

Justification: We have reviewed the NeurIPS Code of Ethnics, and the research conducted in the paper conforms with the NeurIPS Code of Ethnics.

#### Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

### 10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [NA]

[NA]

Justification: The paper is foundational research and not tied to particular applications currently, so there are no potential societal impacts of the work performed.

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.

- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

# 11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: The paper poses no such risks.

#### Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

### 12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: The creators or original owners of assets (e.g., code, data, models), used in the paper, are properly credited, and the license and terms of use are explicitly mentioned and properly respected.

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.

- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

### 13. New assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [No]

Justification: The paper has not released new assets at present.

### Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

# 14. Crowdsourcing and research with human subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: The paper does not involve crowdsourcing nor research with human subjects. Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

# 15. Institutional review board (IRB) approvals or equivalent for research with human subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: The paper does not involve crowdsourcing nor research with human subjects. Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

# 16. Declaration of LLM usage

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or originality of the research, declaration is not required.

Answer: [NA]

Justification: The core method development in this research does not involve LLMs as any important, original, or non-standard components.

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy (https://neurips.cc/Conferences/2025/LLM) for what should or should not be described.