

Seeing Motion Through Polarity for Event-based Action Recognition

Supplementary Material

001 This document provides more details of experimental set-
002 ting, visualization example and prompt setting. The con-
003 tents are meticulously organized in a manner that facilitates
004 ease of understanding and reference, as detailed in the sub-
005 sequent sections :

- 006 • Experimental Setting (cf. §A);
- 007 • Visualization Example (cf. §B).
- 008 • Prompt Setting (cf. §C).

009 A. Experimental Setting

010 We evaluate our research on three benchmarks, including
011 SeAct [5], DVS Action [3] and THU^{E-ACT}-50-CHL [2].

012 **SeAct** dataset is the first semantic-abundant dataset
013 for event-text action recognition. It is collected with a
014 DAVIS346 event camera, whose resolution is 346×260 .
015 It contains 58 actions under four themes, each action is ac-
016 companied by an action caption of less than 30 words gen-
017 erated by GPT-4 to enrich the semantic space of the original
018 action labels.

019 **DVS Action** dataset contains 10 action categories with
020 450 recordings, all captured using the DAVIS346 camera.
021 The action was conducted in an unoccupied office environ-
022 ment, involving the participation of 15 subjects who per-
023 formed a total of 10 different actions.

024 **THU^{E-ACT}-50-CHL** dataset is under challenging
025 scenarios. It contains 58 actions and comprises 2,330
026 recordings featuring 18 students, captured from varied per-
027 spectives utilizing a DAVIS346 event camera in two distinct
028 scenarios: a long corridor and an open hall. Every recording
029 lasts from 2 to 5 seconds and maintains a spatial resolution
030 of 346×260 .

031 For event-to-image reconstruction, we follow EMP [1]
032 to adopt the most representative lightweight network-
033 Firenet [4]. For the event representation, we adopt the
034 method from [5] to reformat the event stream into three-
035 channel event frames. For the text prompts, we utilize a
036 combination of the hand-crafted text prompt and the learn-
037 able text prompt, maintaining consistency with ExACT [5].

038 B. Visualization Example

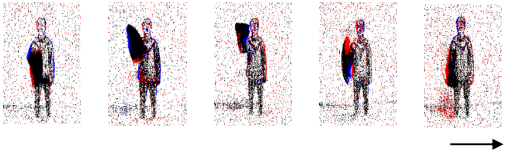
039 **Example for Zero-Shot EAR.** As shown in 1, we illus-
040 trate the complete prompting scheme utilized when test-
041 ing the Multi-modal Large Language Models (MLLMs) for
042 event-based action recognition. The prompt structure is
043 kept consistent across all tested MLLMs, comprising a Sys-
044 tem Prompt and a User Prompt. The MLLM’s output is the
045 predicted label category, which is then compared against the
046 ground-truth label to calculate the accuracy.

Instruction Examples for Zero-Shot EAR

System prompt:
You are an expert analyst in the field of event camera vision, specializing in understanding spatiotemporal dynamics. Your task is to deduce complete human activities by analyzing a series of event frames. **Action belongs to the following categories:** 'arm crossing', 'get-up', 'jumping', 'kicking', 'picking up', 'sit-down', 'throwing', 'turning around', 'walking', 'waving'.

User prompt:
What action is the person in the figure performing, as indicated by the event frames captured by the event camera?

Input:



→ T

Output: 'throwing' ❌

Figure 1. Visualization of instruction examples for zero-shot action recognition.

047 C. Prompt Setting

048 **Complete Progressive Prompting in PMR.** As shown
049 in 2, we illustrate the complete progressive prompting
050 scheme employed in our proposed Polarity Motion Reasoner (PMR). This scheme is designed to explicitly empha-
051 sise that positive polarity represents the forward tendency of
052 motion, while negative polarity signifies the trailing persis-
053 tence of movement. The generation of the final textual de-
054 scription for polarity-enhanced motion is guided by a three-
055 stage progressive instruction pipeline: (1) *Observation* (fo-
056 cusing on raw input cues), (2) *Thinking* (inferring motion
057 characteristics), and (3) *Synthesis and Prediction* (selecting
058 the optimal motion tendency description).
059

060 References

- 061 [1] Meiqi Cao, Xiangbo Shu, Xin Jiang, Rui Yan, Yazhou Yao,
062 and Jinhui Tang. Exploiting frequency dynamics for enhanced
063 multimodal event-based action recognition. In *Proceedings of
064 the IEEE/CVF International Conference on Computer Vision
065 (ICCV)*, pages 5969–5979, 2025. 1
066 [2] Yue Gao, Jiaxuan Lu, Siqi Li, Nan Ma, Shaoyi Du, Yipeng Li,
067 and Qionghai Dai. Action recognition and benchmark using

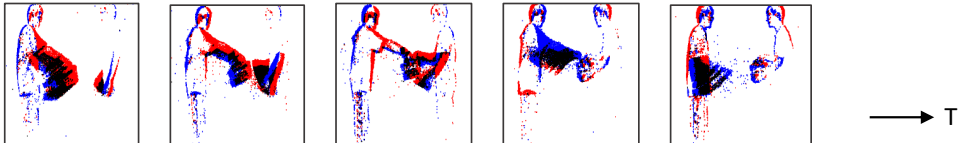
Instruction Examples for PMR	
<p>System prompt:</p> <p>You are an expert analyst in the field of event camera vision, specializing in understanding spatiotemporal dynamics. Your task is to deduce complete human activities by analyzing a series of event frames. In this context, polarities indicate motion: positive polarity (red) events denote the leading edge of movement, and negative polarity (blue) events denote the trailing edge of movement. For each new key frame, you must use the following structured reasoning process to update and expand your ongoing analysis:</p> <p>Step 1: Observation Locate the person in the current frame and identify their status.</p> <p>Step 2: Thinking</p> <p>a) Analyze Instantaneous Motion: Determine the immediate direction and nature of movement (e.g., rotation, swinging, pushing).</p> <p>b) Integrate Temporal Context: Link the just observed instantaneous movement with the broader action narrative built from previous frames.</p> <p>(Example reasoning: "In the previous frame, we observed a dense cluster of red moving upward (arm raising). In the current frame, we see a red arc moving downward. This indicates the motion entering a second phase: a downward swing, directly resulting from the initial arm lift.")</p> <p>Step 3: Synthesis and Prediction</p> <p>Summarize the complete narrative motion (including body areas): Briefly restate the entire event sequence you've observed from the first frame to the current frame. For example: 'First, the object was lifted, remained stationary, and is now swinging forward with the arm.'</p> <p>Adjust your prediction: Based on the comprehensive narrative so far, provide your most plausible prediction of the executed complete action. As more frames appear, this prediction should become more confident or be adjusted.</p>	
<p>User prompt:</p> <p>Please analyze the provided sequence of event frames. Determine the final output format: Your entire response MUST be organized according to the following structure:</p> <p>Thinking (Frame 1→T):</p> <p>Instantaneous Motion: [analysis from Step 2a for the current frame] Temporal Integration: [key analysis from Step 2b, clearly linking current frame with past frames] Overall Motion Narrative: [only describe the critical entire movement with motion regions, formatted as: {Movement Region, Movement Trend}->{Movement Region, Movement Trend}->...] Prediction: [final, concise action derived from the complete narrative] !Check!: As the sequence progresses, your prediction should become more specific, but if evidence is lacking, always choose a general action (e.g., 'swinging') over a possibly incorrect specific action (e.g., 'playing tennis').</p>	
<p>Input:</p>	
<p>Output:</p> <p>{Left Person's Arm, Extends Forward}->{Object, Transfers from Left to Right Person}->{Left Person's Arm, Retracts}</p>	

Figure 2. Visualization of instruction examples for PMR.

068
069
070
071
072

- event cameras. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(12):14081–14097, 2023. 1
- [3] Shu Miao, Guang Chen, Xiangyu Ning, Yang Zi, Kejia Ren, Zhenshan Bing, and Alois Knoll. Neuromorphic vision datasets for pedestrian detection, action recognition, and fall

detection. *Frontiers in Neurobotics*, 13:38, 2019. 1

- [4] Cedric Scheerlinck, Henri Rebecq, Daniel Gehrig, Nick Barnes, Robert Mahony, and Davide Scaramuzza. Fast image reconstruction with an event camera. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer*

073
074
075
076
077

- 078 *Vision*, pages 156–163, 2020. 1
- 079 [5] Jiazhou Zhou, Xu Zheng, Yuanhuiyi Lyu, and Lin Wang. Ex-
080 act: Language-guided conceptual reasoning and uncertainty
081 estimation for event-based action recognition and more. In
082 *Proceedings of the IEEE/CVF Conference on Computer Vi-*
083 *sion and Pattern Recognition*, pages 18633–18643, 2024. 1