

# Seeing without Pixels: Perception from Camera Trajectories

Zihui Xue<sup>1\*, 2</sup> Kristen Grauman<sup>2</sup> Dima Damen<sup>1</sup> Andrew Zisserman<sup>1</sup> Tengda Han<sup>1</sup>  
<sup>1</sup>Google DeepMind <sup>2</sup>The University of Texas at Austin

## Abstract

Can one perceive a video’s content without seeing its pixels, just from the camera trajectory—the path it carves through space? This paper is the first to systematically investigate this seemingly implausible question. Towards this end, we propose a contrastive learning framework to train CamFormer, a dedicated encoder that projects camera pose trajectories into a joint embedding space, aligning them with natural language. We find that, contrary to its apparent simplicity, the camera trajectory is a remarkably informative signal to uncover video content. In other words, “how you move” can indeed reveal “what you are doing” (egocentric) or “observing” (exocentric). We demonstrate the versatility of our learned CamFormer embeddings on a diverse suite of downstream tasks, ranging from cross-modal alignment to classification and temporal analysis. Importantly, our representations are robust across diverse camera pose estimation methods, including both high-fidelity multi-sensored and standard RGB-only estimators. Our findings establish camera trajectory as a lightweight, robust, and versatile modality for perceiving video content.<sup>1</sup>

## 1. Introduction

One sees the environment not with the eyes but with the eyes-in-the-head-on-the-body-resting-on-the-ground.

James J. Gibson

We start with a compelling, perhaps even counter-intuitive, question: can a camera’s trajectory, devoid of all pixels, be informative enough to reveal the content of a video? Consider the challenge in Fig. 1. At first glance, matching these simple, abstract curves to specific human actions seems difficult. The key to this puzzle lies in a fundamental principle: human perception is active. We move to see, turning the visual input into an intentional sub-sampling of our surroundings. This principle extends directly to digital capture. Every video is a human-guided sub-sampling of the world, structured by the creator’s in-

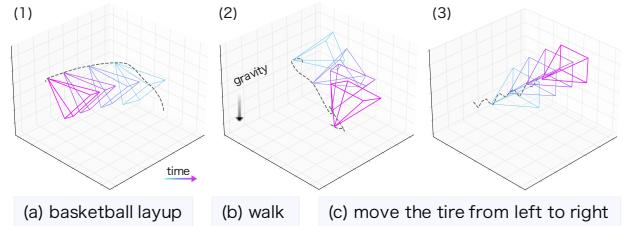


Figure 1. Can you guess which action goes with which camera pose trajectory? In this paper, we find that camera trajectory carries rich information about the video’s content, in both egocentric and exocentric settings. Answers are given in the next page.<sup>2</sup>

tent. This intent is what creates the distinct, recognizable motion signatures in Fig. 1 [spoiler alert]: the upward tilt for a basketball layup, the downward left-to-right sweep for moving a tire, and the rhythmic, forward-moving oscillation for walking, are all potential fingerprints of the semantic action, physically written into the trajectory.

Despite its potential, the intentional signal in camera trajectories has been overlooked. The dominant paradigm learns high-level semantics by training visual encoders on massive-scale video-text data via contrastive learning [3, 38, 51, 68, 69, 78]. While other non-visual modalities like audio [21], IMU [10, 43, 44], thermal, depth images [18] and touch [72] have been explored to replace or complement vision, they often require specialized hardware and thus cannot be obtained retroactively from existing videos. In contrast, camera trajectory, the continuous sequence of camera poses (rotation and translation) over time, is a lightweight and privacy-preserving signal that can be estimated directly from the video itself. We therefore propose to elevate the camera trajectory as a novel modality for video perception, both on its own and as a powerful complement to vision.

The camera trajectory is, of course, not a new concept in computer vision. However, it has been commonly used as a geometric tool, for tasks such as 3D reconstruction and visual odometry [22, 25, 32, 37, 66, 70]. This historical focus has left a fundamental question unexplored: What semantic information, if any, is encoded within the camera trajectory itself? To our knowledge, the camera trajectory has received limited attention as a direct source of evidence for perceiving video content, a gap we address in this paper.

\*Work done during internship at Google DeepMind.

<sup>1</sup>Project webpage: <https://sites.google.com/view/seeing-without-pixels>.

However, learning to “interpret” camera trajectories presents fundamental challenges. The primary prerequisite—access to high-quality data—has long been a blocker. High-quality poses were historically difficult to obtain: hardware-based solutions were not portable or had low sampling rates [28, 77], while conventional estimation methods were computationally expensive or struggled with accuracy [2, 11, 46, 56]. Furthermore, the camera trajectory is a low-dimensional, information-sparse signal, posing a conceptual challenge whether it can be informative to disambiguate the many actions a trajectory might represent.

With regards to the data acquisition challenge, fortunately, recent advances in high-fidelity hardware [15] and accurate pose estimation methods [25, 37, 66, 70] have made large-scale, high-quality camera trajectory data accessible for the first time, creating the conditions for this modality to emerge. With large-scale paired trajectory-text data now at hand, we propose to pre-train a dedicated trajectory encoder, which we term *CamFormer*, to project camera trajectories into a joint embedding space with text, using a contrastive learning framework. Here, the text refers to action narrations or descriptive video captions detailing the video’s content. Next, to address the semantic ambiguity inherent in camera trajectories, we propose a contextualized trajectory encoding that incorporates extended temporal context for disambiguation.

We structure our investigation to cover two distinct scenarios. We first analyze the *egocentric (first-person)* setting, where a wearable camera’s trajectory offers a direct correlation with the recorder’s action. We then analyze *exocentric (third-person)* videos, where the trajectory is decoupled from the actor, and reflects the recorder’s attention as an observer. To investigate both domains, our evaluation is structured around three core capabilities (cross-modal alignment, classification, and temporal analysis) and spans multiple dimensions, including semantic granularity (coarse activities, fine keysteps) and task types (retrieval, classification, localization). Across this comprehensive suite of 10 tasks on 5 datasets, CamFormer effectively unlocks the semantic potential of the camera trajectory, delivering consistent gains ranging from +3.2% to +13.2%. The trajectory proves to be a powerful standalone signal: our lightweight CamFormer is capable of outperforming computationally-heavy vision models in key scenarios. It also excels as a valuable complementary signal, providing the best overall performance when fused with vision. Finally, we show that CamFormer is robust across various camera pose sources, from high-fidelity multi-sensor SLAM to standard video-only estimates, demonstrating its practical utility.

## 2. Related Work

**Multimodal Contrastive Learning.** Our world is inherently multimodal. The success of vision-language mod-

els like CLIP [52] establish a de facto standard for binding our rich visual world to the semantic structure of text [38, 51, 68, 69], using a contrastive objective [48] that aligns corresponding pairs in a shared embedding space. The same principle has been applied to a broader set of modalities, seeking to connect audio [21], IMU [10, 43, 44], thermal, depth images [18] and touch [72] to a shared semantic space. Conspicuously absent from this list, however, is the camera pose trajectory, despite being an intrinsic property of any video recording. Our work addresses this gap, introducing it as a new modality and demonstrating its immense potential for semantic representation learning through a series of downstream tasks.

**Action Understanding with Egocentric Motion.** Motion has long been recognized as a helpful signal in egocentric videos. A few works [1, 47, 54, 60] have shown that egocentric motion representations, such as optical flow [29, 30, 36, 59, 65], can aid action recognition. Another line of work [13, 41, 49, 62, 63, 76] captures motion signals explicitly using IMU sensors mounted on the camera wearer’s head or limbs. Beyond action recognition, egocentric motion has also been shown to correlate with other semantic properties, such as physical forces [50] and camera-wearer identity [74]. While these works support the premise that egocentric motion is informative, they are confined to the egocentric domain and a narrow set of tasks. Our investigation takes a broader view by systematically investigating the camera trajectory, across both egocentric and exocentric domains, and on a diverse suite of downstream tasks.

**Camera Pose in 3D Vision.** Camera pose estimation is a fundamental task of 3D vision. Recent advances deliver increasingly precise trajectories, driven by both multi-sensor hardware like Meta Aria glasses [15] and learning-based models [25, 31, 37, 66, 70, 71] that infer pose from monocular videos. Our work is orthogonal to these efforts: we focus on interpreting the resulting trajectories, with performance naturally benefiting as estimation continues to improve. On the application side, camera pose has powered tasks such as novel view synthesis [34], 3D reconstruction and mapping [27], or serving as a conditional prior to guide video generation [17, 24, 32, 35]. While works on 3D human motion estimation [26, 33, 40, 73, 75] and hand forecasting [23] share a similar premise—that camera motion is a signature of the actor’s movement—their goal remains the estimation of the physical body / hands. Our work is the first to broaden this perspective, proposing that the trajectory encodes rich semantic signals for video perception.

**Generating Camera Motion Descriptions.** Recent work [16, 39, 67] trains large multimodal models (LMMs) to generate textual descriptions of camera cinematography from video input (e.g., “zoom”, “pan”). While these works

---

<sup>2</sup>Answer: (1)-(c), (2)-(a), (3)-(q)

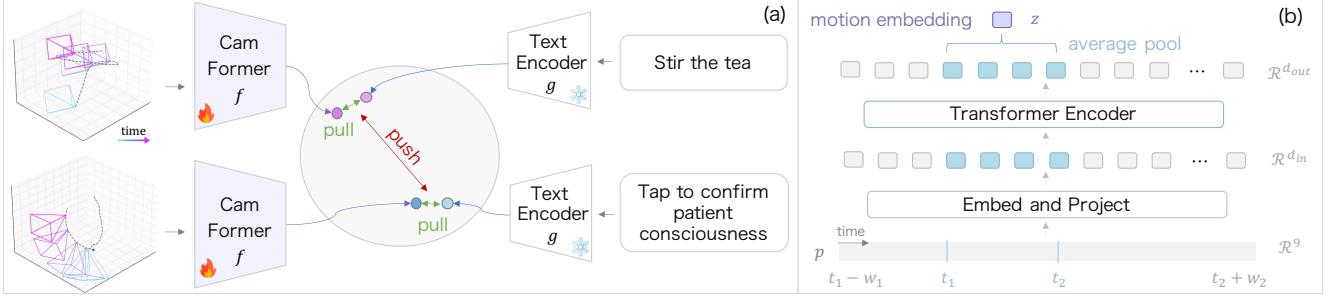


Figure 2. **Unlocking Semantic Information Hidden in Camera Trajectories.** (a) We propose contrastive pre-training on paired (trajectory, text) data. Our model, CamFormer, is trained to map camera trajectories into a joint semantic space, aligning them with natural language. (b) We propose contextualized trajectory encoding that incorporates extended temporal context to disambiguate the local action.

treat camera motion as a video attribute to describe, we treat the trajectory as a semantic signal to interpret, which our experiments prove better decodes video content.

### 3. Method

This work addresses the untapped problem of how to unlock the semantic information within camera pose trajectories. Our core methodology is to learn a dedicated encoder to map camera trajectories into a shared semantic space with text (Sec. 3.1). We then apply the learned camera trajectory embeddings to a suite of downstream tasks to assess their utility for perceiving video content (Sec. 3.2).

#### 3.1. Learning Camera Trajectory Embeddings

We apply multimodal contrastive learning [52] to a new modality pairing: the camera trajectory of a video clip and the text description of that clip’s content.

**Problem Formulation.** Formally, let  $\mathbf{v}$  be a video clip spanning the time interval  $[t_1, t_2]$ , and  $\mathbf{t}$  be its paired text description. Let  $\mathbf{p} \in \mathbb{R}^{N \times 9}$  denote the corresponding camera pose trajectory derived from  $\mathbf{v}$ , where  $N$  is the number of pose samples extracted from the clip’s duration  $t_2 - t_1$  at a given sampling rate  $s$ , i.e.,  $N = (t_2 - t_1) \times s$ . Each pose is represented as a 9D relative vector (3D translation and 6D continuous rotation representation [79]), computed with respect to the sequence midpoint. The representation choice is justified in our ablation (cf. Supplementary).

**Training Objective.** We employ InfoNCE [48] loss, where the matching pairs of camera trajectory and text in a batch are treated as positives and all other pairwise combinations in the batch are regarded as negatives. For a batch of  $B$  examples,  $\{(\mathbf{p}_i, \mathbf{t}_i)\}_{i=1}^B$ , we optimize the loss  $\mathcal{L} = \mathcal{L}_{P \rightarrow T} + \mathcal{L}_{T \rightarrow P}$ :

$$\mathcal{L}_{P \rightarrow T} = -\frac{1}{B} \sum_{i=1}^B \log \frac{\exp(f(\mathbf{p}_i) \cdot g(\mathbf{t}_i)/\tau)}{\sum_{j=1}^B \exp(f(\mathbf{p}_i) \cdot g(\mathbf{t}_j)/\tau)}$$

where  $\mathcal{L}_{T \rightarrow P}$  is the symmetric text-to-trajectory loss, and  $\tau$  is the temperature hyperparameter. The text encoder  $g$  is

the pre-trained and frozen encoder from CLIP [52], which provides a fixed semantic target and allows  $f$  to learn representations grounded in its robust embedding space.

**Model Architecture.** Our camera trajectory encoder  $f$  (CamFormer), is a light-weight Transformer [64]. The input 9D pose sequence  $\mathbf{p} \in \mathbb{R}^{N \times 9}$  is first linearly projected into  $d_{in}$ -dimensional embeddings  $\mathbb{R}^{N \times d_{in}}$ . After adding positional embeddings, the full sequence is processed by a series of Transformer blocks to fuse temporal information. Finally, the output features are temporally mean-pooled to a vector in  $\mathbb{R}^{d_{in}}$ , and projected by a linear layer to  $d_{out}$  dimensions to produce the final vector  $\mathbf{z} = f(\mathbf{p}) \in \mathbb{R}^{d_{out}}$ .

**Contextualized Trajectory Encoding.** A unique challenge in encoding camera trajectories is the low information density. Comparing with a short video clip (*e.g.* 1 second), a camera trajectory with a similar duration carries much sparser information and could be semantically ambiguous. We address this challenge with a contextualized trajectory encoding, which extends the sequence length beyond the immediate temporal window to capture broader temporal context, thereby disambiguating the central action.

Formally, we extend the base window  $[t_1, t_2]$  by a total duration  $w$ .  $w$  is then randomly split into  $w_1$  and  $w_2$  (such that  $w_1 + w_2 = w$ ), resulting in an extended, temporally-shifted window  $[t_1 - w_1, t_2 + w_2]$ . The trajectory in this extended window serves as the final input to our CamFormer  $f$ . After  $f$  processes this entire sequence, the final embedding  $\mathbf{z}$  is produced by mean-pooling only the  $N$  output features corresponding to the original window  $[t_1, t_2]$ . This strategy infuses the local representation with global context without diluting it with potentially irrelevant adjacent actions. The encoding process is shown in Fig. 2 (b).

#### 3.2. Applications of Camera Trajectory Embeddings

After pre-training CamFormer  $f$  on paired camera trajectories and text, we test the power of its representation across a diverse suite of downstream tasks:

- **Cross-modal Alignment.** Through text retrieval, we directly test the learned alignment between a camera trajec-

Table 1. Datasets used in our analysis. Black checkmarks ( $\checkmark$ ) denote data provided by the original dataset, while green checkmarks ( $\checkmark$ ) reflect our enhancement (estimated camera trajectories). The subscript  $\times n$  specifies the number of available trajectories sources. We provide three additional pose sources (MegaSaM [37], ViPE [25] and  $\pi^3$  [70]) to enrich Ego-Exo4D and UCF101, and one ( $\pi^3$ ) to FineGym. “label” denotes downstream task annotations (e.g., action labels). The green-highlighted rows denote our pre-training datasets.

Dataset	Domain	Trajectory	Text	# Hrs
Ego-Exo4D [20]	ego	$\checkmark \checkmark \times 3$	narration & label	221.3
Nymeria [42]	ego	$\checkmark$	narration	38.6
DynPose-100K [53]	exo	$\checkmark \times 2$	caption	157.5
UCF101 [61]	exo	$\checkmark \times 3$	label	27.0
FineGym [57]	exo	$\checkmark \times 1$	label	92.8

tory and the video’s content, as represented by its paired text (i.e., action narrations or video captions).

- **Downstream Classification.** We assess our embeddings on a spectrum of classification tasks. This includes identifying coarse-grained scenarios reflecting the overall activity (e.g., distinguishing between cooking, dancing, or basketball) and fine-grained keysteps (e.g., within a cooking activity, identifying cutting, pouring water and washing fruit). Additionally, we evaluate the ability to discern motion signatures indicative of performer skill levels via proficiency estimation (i.e., beginner or expert).
- **Temporal Analysis.** We examine tasks that require precise temporal information, including temporal action localization and the recognition of periodic patterns for repetitive action counting.

The pre-trained CamFormer is applied to these diverse tasks using two standard evaluation paradigms: (1) as a frozen feature extractor (e.g., linear probing); and (2) via end-to-end fine-tuning, where CamFormer is trained jointly with a linear head. See Supp. for details.

## 4. Experimental Setup

**Bridging the Data Gap.** We observe a trajectory-semantic data gap: traditional 2D video understanding datasets are vast and semantically rich but lack camera pose trajectories [9, 19, 57, 61], while 3D datasets that include poses often focus on static scenes and navigation tasks [5, 8, 55], rather than depicting diverse human actions or high-level semantics. Fortunately, recent advances in both hardware capture and software-based pose estimation [15, 25, 37, 66, 70] allow us to source and assemble these components. Table 1 summarizes our data effort to bridge the gap.

As shown in the table, we adopt two large-scale datasets for pre-training. For the egocentric domain, we adopt the large-scale Ego-Exo4D [20], which provides time-stamped narrations and high-quality dense camera trajectories from Meta Aria glasses [15] (specifically, visual-inertial pose

estimates from the device’s SLAM cameras). As these are head-mounted, the camera trajectory serves as a direct and high-fidelity proxy for the wearer’s head motion. We only use its egocentric videos, as the dataset’s exocentric cameras are static and thus uninformative for our motion-based analysis. For the exocentric domain, we pre-train on DynPose-100K [53] data, where we utilize video captions from Panda70M [6] and the two available pose sources (the original one provided by the dataset and an alternate one estimated by ViPE [25]). Our downstream evaluation suite includes both Ego-Exo4D and DynPose-100K, supplemented by Nymeria [42] (another egocentric dataset collected by Aria glasses for full-body motion understanding) and two standard exocentric action recognition datasets (UCF101 [61] and FineGym [57]).

We enrich these datasets with estimated camera trajectories. On Ego-Exo4D, we supplement the dataset’s original camera poses with trajectories from three leading pose estimators (MegaSaM [37], ViPE [25], and  $\pi^3$  [70]). This multiple-source testbed not only allows us to benchmark model robustness across pose sources, but also offers a novel, semantic-based evaluation for the pose estimators themselves (i.e., which estimator yields the best downstream semantic performance, cf. Table 3). We also apply this estimation pipeline to UCF101 and FineGym to generate the pose trajectories that they lack, thereby enabling our analysis on these standard action benchmarks.

**Implementation.** We sample input camera poses at 5-30 Hz, subject to datasets. For our CamFormer, we set the internal projection dimension  $d_{in}$  to 128 and the final output dimension  $d_{out}$  to 512, which matches our frozen text CLIP encoder. The transformer encoder consists of 4 layers, each with 4 attention heads, a 256-dimensional feed-forward network, and a dropout rate of 0.1. During training, the context duration  $w$  is uniformly sampled from  $\mathcal{U}(0, w_{\max})$ , where  $w_{\max} = 8s$  is the maximum context length. During inference, we test various values of  $w$  to analyze the impact of temporal context (cf. Sec. 5.3).

## 5. Results

We structure our findings as a series of questions, exploring the egocentric (Sec. 5.1) and exocentric (Sec. 5.2) domains, followed by further analysis (Sec. 5.3). Fig. 3 (a) provides a high-level results overview, comparing CamFormer against base methods / models across 10 downstream tasks, and highlighting its consistent performance gains that we will detail throughout this section.

### 5.1. CamFormer for Egocentric Videos

Our analysis in the egocentric domain reveals a strong, direct link between the camera’s trajectory and the recorder’s activity, which our CamFormer effectively captures. We see this in Fig. 3 (b), which visualizes PCA embeddings

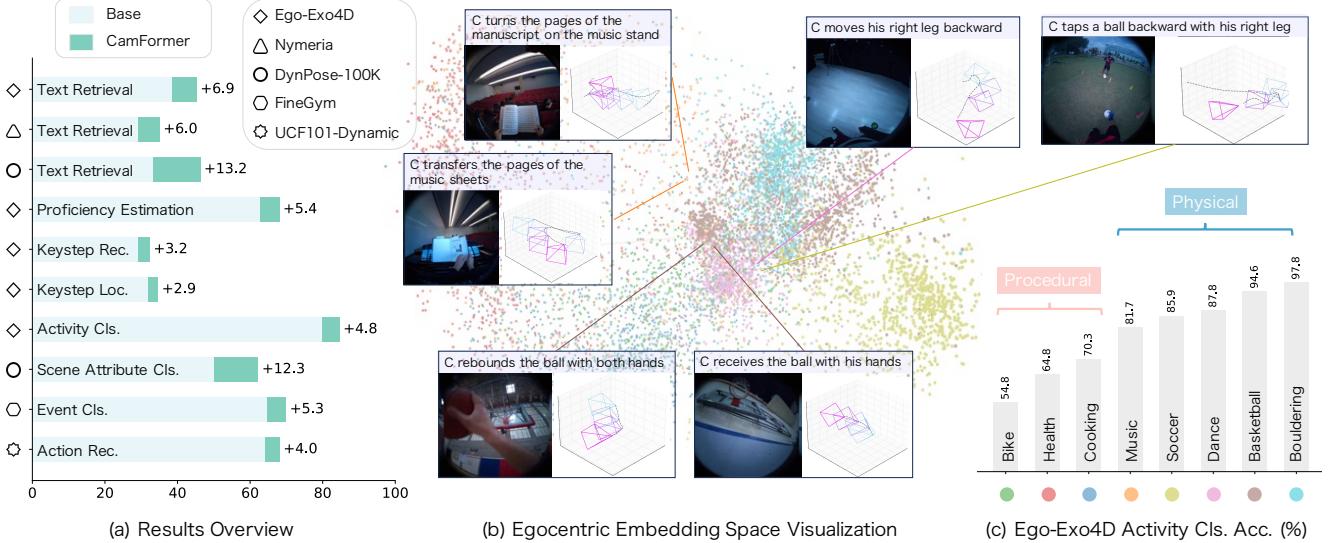


Figure 3. (a) Quantitative Results Overview: we summarize CamFormer’s performance against base methods / models on 10 downstream tasks across 5 datasets, demonstrating its consistent performance advantages; (b) A PCA visualization of CamFormer embeddings on unseen Ego-Exo4D trajectories, colored by the dataset’s 8 activity labels. (Note: CamFormer only takes the trajectory as input; video clips and text are shown for interpretation only); (c) Per-class activity classification accuracy plot on Ego-Exo4D reveals a performance dichotomy: CamFormer excels at physical activities but is less effective on procedural ones with more subtle camera motions.

of CamFormer’s features on the unseen Ego-Exo4D validation set. The embeddings demonstrate some natural clustering by the dataset’s activity labels (a task for which our model was not explicitly trained). Furthermore, the visualized data pairs show that trajectories close in the embedding space share action semantics, even if their specific motion patterns or overall activity labels are different. This demonstrates that CamFormer learns to decipher the underlying semantic meaning encoded within a camera trajectory through our pre-training.

#### Q: When is camera trajectory most informative?

To quantify the emergent clustering observed in the embedding space, we fine-tune CamFormer for the downstream activity classification task on Ego-Exo4D using its 8 activity labels. The per-class accuracy shown in Fig. 3 (c) reveals a clear pattern: our model excels at recognizing dynamic physical activities, achieving accuracies over 90% for basketball and bouldering, but gets confused among the three procedural activities (i.e., bike repair, health and cooking), where camera motions are more localized and subtle. This observation forms the basis of our analysis below.

#### Q: How does camera trajectory compare to other modalities?

**Evaluation Setup** To serve as our main evaluation for cross-modal alignment, we create a new, challenging text retrieval benchmark from the Ego-Exo4D validation split, comprising 7079 camera trajectory queries. To reduce retrieval ambiguity, we adapt the 5-way multiple-choice ques-

tion (MCQ) format from [38], where the model must select the correct textual description from five candidates. The queries are carefully balanced across the dataset’s 8 activity labels and its “in-view” (iv) / “out-of-view” (ov) visibility annotations to facilitate more detailed analysis. Fig. 4 (up) provides an example of this task, showing an “ovv” case where the action is not directly visible in the egocentric video (see Supp. for more).

Our primary comparison is with established contrastive multimodal methods incorporating common modalities aligned with text: CLAP (audio-text) [14], PRIMUS (IMU-text) [10], CLIP (image-text) [52] and EgoVLPv2 (video-text) [51]. We report on two EgoVLPv2 checkpoints (pre-trained on Ego4D [19] and Ego-Exo4D, respectively) and note that PRIMUS is also pre-trained on Ego-Exo4D, aligning with our setup. As text retrieval is formulated as MCQ, we also report performance of one representative LMM, Gemini-2.5-Pro [7], to highlight the challenging nature of our MCQ task. We stress that this is not a direct comparison due to Gemini’s prohibitive cost and different paradigm.

**Results** Table 2 benchmarks CamFormer’s performance on our 5-way MCQ task. Remarkably, on physical activities, it outperforms computationally-heavy video baselines with a much lower computational cost, proving its strength as a lightweight standalone signal. On procedural activities, CamFormer embeddings provide the best results overall when fused with video (achieved by averaging our camera trajectory features with the best-performing video features, EgoVLPv2 [51]). This confirms the camera trajec-

Table 2. Egocentric Text Retrieval Results. We report top-1 accuracy (%) for a 5-way MCQ on Ego-Exo4D and Nymeria (zero-shot). Chance is 20%. Best results are bolded and second best results are underlined. For Ego-Exo4D, columns show performance on splits where the narration is in-view (iv) or out-of-view (oov) from the *egocentric* perspective, across both physical and procedural activities. For Nymeria, columns show performance on four motion annotation types (legs/feet, focus attention, body posture, hands/arms). The grayed-out Gemini-2.5-Pro row is included to highlight the challenging nature of our benchmark, not as a direct, apples-to-apples comparison. Key takeaways are: (1) As a standalone signal, the camera trajectory is a strong, low-cost modality that outperforms video on physical activities and in visually challenging scenarios (*e.g.*, “oov” and “legs” split); (2) As a complementary signal, for cases where vision is strong (*e.g.*, procedural activities), fusing it with video features (\*) improves over the video-only baseline, proving its non-redundant value.

Method	Modality	# MACs (G)	# Params (M)	Ego-Exo4D [20]						Nymeria [42]							
				Physical		Procedural		all	legs		focus		body		hands		all
				iv	oov	iv	oov		legs	focus	body	hands	body	hands	body	hands	
CLAP [14]	audio	6.80	32.8	21.4	20.1	29.5	33.3	24.6	-	-	-	-	-	-	-	-	
PRIMUS [10]	IMU	0.03	1.4	18.5	25.3	24.7	22.0	23.2	-	-	-	-	-	-	-	-	
CLIP [52]	image	2.95	59.0	25.2	18.2	26.8	21.9	22.9	22.4	22.4	40.5	29.5	28.7				
EgoVLPv2 [51] (Ego4D)	video	89.49	150.7	30.9	24.8	48.7	40.9	34.4	22.1	17.4	45.7	31.7	29.2				
EgoVLPv2 [51] (Ego-Exo4D)	video	89.49	150.7	39.1	25.6	50.5	45.4	38.4	22.8	20.5	39.6	28.5	27.9				
Gemini-2.5-Pro [7]	video	-	-	53.9	29.0	67.4	50.9	48.7	32.6	36.6	60.7	52.7	45.6				
CamFormer embeddings	trajectory	0.02	0.3	<b>56.1</b>	<b>46.4</b>	34.3	32.7	<u>44.8</u>	<b>30.8</b>	<b>33.2</b>	36.3	26.0	<u>31.6</u>				
CamFormer embeddings*	video+trajectory	89.51	151.0	<b>56.0</b>	<b>45.8</b>	<b>51.4</b>	<b>45.9</b>	<b>46.0</b>	<b>30.1</b>	<b>30.6</b>	<b>47.8</b>	<b>34.0</b>	<b>35.6</b>				

tory is a versatile signal, effective both on its own and as a complementary source of information for video. Fig. 4 (up) provides an “oov” example that illustrates its advantage. In the bouldering activity, distinguishing between fine-grained actions like “rise towards wall” and “land on the mat” is difficult from visual frames alone. However, the camera pose trajectory provides an unambiguous downward signal, allowing our model to correctly identify “land on the mat.” This confirms camera trajectory’s unique value in disambiguating actions when visual cues are subtle or misleading.

### Q: How generalizable is the learned CamFormer?

**Evaluation Setup** After CamFormer is pre-trained exclusively on Ego-Exo4D [20], we test this model by directly applying it on Nymeria [42] dataset in a zero-shot manner. We use the entire Nymeria dataset as a test set and create 4000 questions (1000 for each of its four narration types: legs/feet, focus attention, body posture, hands/arms), following the same MCQ design used for Ego-Exo4D.

**Results** Table 2 right columns show that our model demonstrates strong zero-shot generalization to Nymeria, performing effectively across its diverse annotation types. Relative to baselines, it is particularly strong on narrations describing “legs” and “focus attention”—the categories often being out-of-view, where baseline visual models perform near random guess. This confirms our model’s unique ability to understand non-visible actions from their motion signature, echoing our “oov” analysis on Ego-Exo4D.

### Q: How does camera trajectory perform on physical vs. procedural tasks?

**Evaluation Setup** Based on our activity classification findings (Fig. 3 (c)), we perform a deeper analysis on the fol-



Figure 4. Qualitative Text Retrieval Results on egocentric Ego-Exo4D (up) and exocentric DynPose-100K (bottom). Up: A clear downward pose trajectory disambiguates the action of landing, where visual cues are subtle. Bottom: A circling trajectory, common for capturing a scene overview, is correctly associated with the high-level scene description. See Supp. for more qualitatives.

**Table 3. Analysis of Model Robustness.** Our CamFormer, despite being pre-trained only on Aria poses, generalizes effectively to process various estimated poses. Left (Ego-Exo4D activity classification): Models initialized from our CamFormer consistently outperform their counterparts trained from scratch. Right (Ego-Exo4D keystep recognition): Fusing camera trajectory features with video consistently outperforms the video-only baseline (29.17%), with performance gains shown in brackets. This demonstrates the complementary value of camera trajectory.

Pose Source	Activity Cls.			Keystep Rec.	
	Pretrain ✗	Pretrain ✓	Trajectory	Video+Trajectory	
MegaSaM [37]	53.67	60.83 (+7.16)	11.49	32.60 (+3.43)	
ViPE [25]	60.83	66.15 (+5.32)	12.26	32.21 (+3.04)	
$\pi^3$ [70]	61.47	66.15 (+4.68)	12.76	32.83 (+3.66)	
Aria [15]	61.83	71.28 (+9.45)	14.07	32.37 (+3.20)	

lowing targeted Ego-Exo4D tasks: (1) proficiency estimation on binary labels of expert vs. non-expert for the two physical activities (bouldering, dancing); (2) keystep recog-

Table 4. Exocentric Text Retrieval Results on DynPose-100K. We report top-1 accuracy (%) for a 5-way MCQ. Our CamFormer not only performs well above random chance (20%) across both pose sources (original [53] and ViPE [25]), confirming a meaningful trajectory-semantic link in the exocentric domain, but also surpasses the “camera description” baselines, demonstrating the effectiveness of interpreting the camera trajectory directly.

Method	Modality	Acc. (%)
Qwen-2.5-VL-7B on CameraBench [16]	cam desc. (text)	27.8
ShotVL-7B [39]	cam desc. (text)	33.1
CamFormer embeddings (original)	trajectory	36.2
CamFormer embeddings (ViPE)	trajectory	<b>46.3</b>

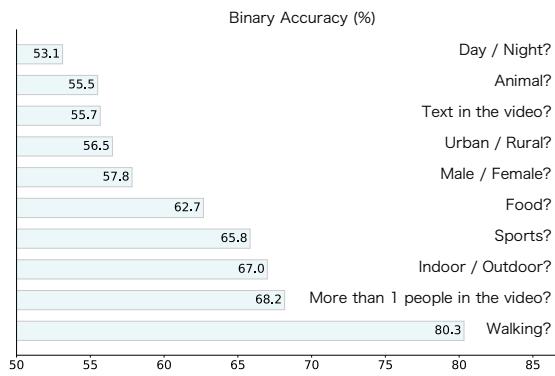


Figure 5. Scene Attribute Classification Results on DynPose-100K. The results reveal a clear spectrum of what can and cannot be inferred from an observer’s camera trajectory.

nition on the 278 fine-grained keystep labels, which are specific to the 3 procedural activities (bike repair, health, cooking); and (3) temporal action localization, using keystep boundaries from the same procedural activities.

**Results** This analysis confirms the pattern from Fig. 3 (c), revealing a two-fold role for the camera trajectory:

- For *physical activities*, camera trajectory is a powerful standalone signal. On proficiency estimation task, CamFormer achieves an average +5.4% gain over the strong video baseline [4] (Fig. 3 (a)), effectively capturing the motion signatures of expertise.
- For *procedural activities*, where motion patterns can be ambiguous, camera trajectory provides valuable complementary information. Fusing CamFormer features with the strong EgoVLPv2 [51] video features on keystep tasks (Fig. 3 (a)) yields a +3.2% accuracy gain in recognition and a +2.9 mIoU@0.3 gain in localization.

Full results for these tasks are available in Supplementary.

#### Q: How robust is CamFormer to estimated poses?

**Evaluation Setup** Up to now CamFormer is pre-trained on the given Aria camera poses [15], which are highly accurate SLAM-based multi-sensor estimates. We further investigate the robustness of CamFormer on more accessible, video-based camera pose estimates. We select two representative tasks on Ego-Exo4D for this analysis: coarse-grained

activity classification, which we evaluate with end-to-end fine-tuning, and fine-grained keystep recognition, which we evaluate using a linear SVM on frozen features.

**Results** Table 3 shows that our CamFormer, pre-trained only on Aria trajectories, is robust and generalizes well across three different estimated poses. First, for activity classification, the benefit of our pre-training holds: a model initialized with our Aria-pretrained CamFormer remains superior to one trained from scratch, even when evaluated on these video-based camera estimates. Second, for keystep recognition, the trajectory’s complementary power persists, as the video-trajectory fusion approach (using estimated poses) continues to outperform the video-only baseline.

Among the estimators,  $\pi^3$  [70] emerges as the top performer. Yet, there are still gaps of these video-only pose estimation methods to Aria trajectories, which highlight the need for further improvements in estimation accuracy. Our analysis also introduces a valuable new direction: a novel semantic benchmark that evaluates estimators on their downstream semantic utility, offering a more holistic assessment than geometric error alone.

## 5.2. CamFormer for Exocentric Videos

#### Q: Does the trajectory-semantic link exist in exocentric videos (third-person view)?

Building on our egocentric findings, we now analyze the exocentric domain, where we posit the attentional imprint of the observer is a rich signal for the perceived event. For instance, the quick, reactive pans of an observer filming a soccer match are likely to differ from the slow, steady hold they would use to film a museum painting.

**Evaluation Setup** We evaluate the exocentric CamFormer (trained on DynPose-100K) on four tasks: (1) text retrieval, where we use the same 5-way MCQ setup as in the egocentric domain, on a 1000-query DynPose-100K test set. For this task, we compare against a two-stage “camera description” baseline, where we first prompt a specialized LMM [16, 39] to generate camera descriptions from the video, then adopt a second LLM to answer the MCQ using that text. (2) scene attribute classification, where we design a list of 10 attribute questions and automatically label a held-out set on DynPose-100K via LMM prompting; (3) event classification on FineGym using its 4 event labels; and (4) action recognition on UCF-Dynamic, our custom 8-class subset of UCF101. See Supp. for details.

#### Results

- *Text Retrieval.* Our exocentric results (Table 4) confirm that a meaningful link between observer trajectories and video content exists. Our CamFormer performs well above random chance (20%) across both pose sources (ViPE and the original). Importantly, it outperforms the multi-stage camera description baselines. This demon-

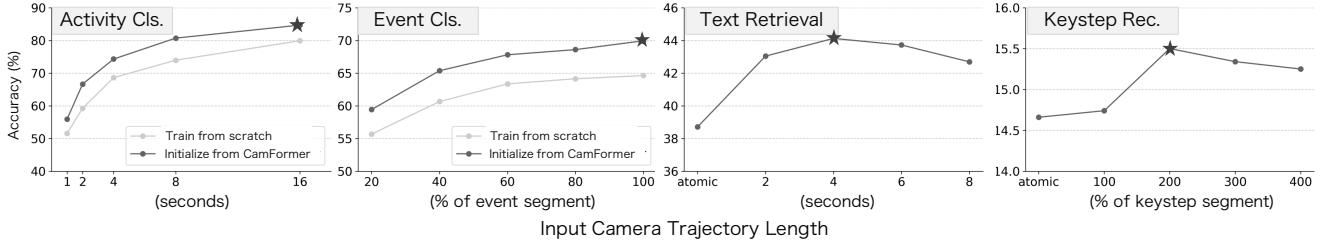


Figure 6. Analysis of Contextualized Trajectory Encoding. For the global label tasks (left two plots), performance monotonically increases with longer input sequences, as more evidence helps disambiguate the activity / event. For the localized label tasks (right two plots), performance peaks at an optimal context length, and then declines as excessive, irrelevant motion acts as noise.

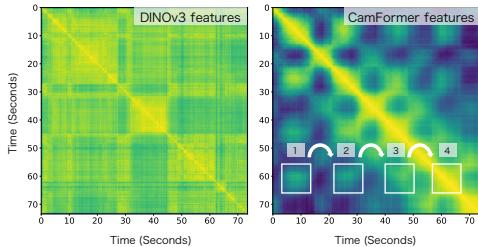


Figure 7. Comparison of temporal self-similarity maps from DINOv3 [58] framewise features (left) and our CamFormer features (right) on an egocentric cutting sequence from Ego-Exo4D. Our map reveals clear, periodic diagonal streaks that correspond to the number of cuts. The visual feature map produced by DINOv3, in contrast, is unstructured and fails to capture the repeating action.

- strates the effectiveness of interpreting camera trajectories directly and suggests that current LMMs have yet to perfect the generation of accurate camera motion descriptions. Fig. 4 (bottom) shows one qualitative example, where our CamFormer correctly links a circling overview motion to its corresponding scene description.
- **Scene Attribute Classification.** Echoing the activity-level analysis in the egocentric domain, we test our exocentric CamFormer’s ability to classify a range of binary scene attributes on DynPose-100K (Fig. 5). The analysis reveals what can be inferred from an observer’s camera motion. At one end, it provides no discernible signal for static attributes like day vs. night, where performance is near random. At the other end, it is a strong predictor for physical attributes, *e.g.*, achieving over 80% accuracy for classifying if there is people walking in the video.
  - **Event & Action Classification.** As summarized in Fig. 3 (a), our CamFormer learned on DynPose-100K, effectively delivers benefits to semantic classification tasks. Compared to a base model trained from scratch, initializing from CamFormer achieves a +5.3% gain on FineGym event classification and a +4.0% gain on UCF-Dynamic action recognition. See Supp. for full results.

### 5.3. Further Analysis

**Analysis of Temporal Context Length.** We validate our contextualized trajectory encoding strategy (Fig. 6) by vary-

ing the input camera trajectory length across two settings:

- For *tasks with one global label* (Ego-Exo4D activity and FineGym event classification), the left two plots show that performance steadily improves with longer sequences, as more evidence helps disambiguate the overall activity / event. Furthermore, our model initialized from the pre-trained CamFormer consistently outperforms the train-from-scratch baseline across all sequence lengths, showcasing the benefits of our proposed pre-training.
- For *tasks with temporally localized labels* (Ego-Exo4D text retrieval and keystep recognition), we investigate the impact of including temporal context from outside the labeled time window. The right two plots show that performance initially improves as more context is added, demonstrating that surrounding motion is crucial for understanding an action. This trend reverses when the context window becomes excessively large, as irrelevant motion begins to act as noise. This reveals an optimal context “sweet spot” for these localized tasks.

**Repetitive Action Counting** is an interesting emergent capability of CamFormer. To do this, we compute the self-similarity map of CamFormer’s temporal features (*i.e.*, the output token sequence before the final average pooling). As shown in Fig. 7, our map (right) reveals clear, periodic diagonal streaks for an egocentric cutting sequence, in contrast to the unstructured map from DINOv3 [58] visual features (left). These periodic patterns are precisely the type of fine-grained temporal signal required for counting [12]. See Supp. for more qualitatives and an animation.

## 6. Conclusion

This work challenges the traditional, geometric view of camera trajectory, showcasing its potential as a semantic signal. Our results with CamFormer reveal the trajectory-semantic link exists in both egocentric and exocentric domains. Crucially, CamFormer is robust across diverse pose estimation methods, demonstrating its practical utility for real-world videos. We hope these findings inspire the community to take a new look at the camera trajectory and further explore this promising modality.

## Acknowledgment

We thank Junyu Xie for assistance and valuable suggestions on the camera pose estimation pipeline; Jean-Baptiste Alayrac, Carl Doersch and Ignacio Rocco for constructive feedback; and Ang Cao, Yue Zhao and Lorenzo Torresani for helpful discussions.

## References

- [1] Girmaw Abebe, Andrea Cavallaro, and Xavier Parra. Robust multi-dimensional motion features for first-person vision activity recognition. *Computer Vision and Image Understanding*, 149, 2016. [2](#)
- [2] Sameer Agarwal, Yasutaka Furukawa, Noah Snavely, Ian Simon, Brian Curless, Steven M Seitz, and Richard Szeliski. Building rome in a day. *Communications of the ACM*, 54(10), 2011. [2](#)
- [3] Kumar Ashutosh, Rohit Girdhar, Lorenzo Torresani, and Kristen Grauman. HierVL: Learning hierarchical video-language embeddings. In *CVPR*, 2023. [1](#)
- [4] Gedas Bertasius, Heng Wang, and Lorenzo Torresani. Is space-time attention all you need for video understanding? In *Icmi*, 2021. [7, 3](#)
- [5] Angel Chang, Angela Dai, Thomas Funkhouser, Maciej Halber, Matthias Niessner, Manolis Savva, Shuran Song, Andy Zeng, and Yinda Zhang. Matterport3D: Learning from rgb-d data in indoor environments. *arXiv preprint arXiv:1709.06158*, 2017. [4](#)
- [6] Tsai-Shien Chen, Aliaksandr Siarohin, Willi Menapace, Ekaterina Deyneka, Hsiang-wei Chao, Byung Eun Jeon, Yuwei Fang, Hsin-Ying Lee, Jian Ren, Ming-Hsuan Yang, et al. Panda-70m: Captioning 70m videos with multiple cross-modality teachers. In *CVPR*, 2024. [4, 2](#)
- [7] Gheorghe Comanici, Eric Bieber, Mike Schaeckermann, Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Marcel Blissein, Ori Ram, Dan Zhang, Evan Rosen, et al. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. *arXiv preprint arXiv:2507.06261*, 2025. [5, 6, 1, 3](#)
- [8] Angela Dai, Angel X Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Nießner. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In *CVPR*, 2017. [4](#)
- [9] Dima Damen, Hazel Doughty, Giovanni Maria Farinella, Sanja Fidler, Antonino Furnari, Evangelos Kazakos, Davide Moltisanti, Jonathan Munro, Toby Perrett, Will Price, et al. Scaling egocentric vision: The epic-kitchens dataset. In *ECCV*, 2018. [4](#)
- [10] Arnav M Das, Chi Ian Tang, Fahim Kawsar, and Mohammad Malekzadeh. Primus: Pretraining imu encoders with multi-modal self-supervision. In *ICASSP*. IEEE, 2025. [1, 2, 5, 6](#)
- [11] Andrew J Davison, Ian D Reid, Nicholas D Molton, and Olivier Stasse. Monoslam: Real-time single camera slam. *IEEE transactions on pattern analysis and machine intelligence*, 29(6), 2007. [2](#)
- [12] Debidatta Dwibedi, Yusuf Aytar, Jonathan Tompson, Pierre Sermanet, and Andrew Zisserman. Counting out time: Class agnostic video repetition counting in the wild. In *CVPR*, 2020. [8](#)
- [13] Kiana Ehsani, Daniel Gordon, Thomas Nguyen, Roozbeh Mottaghi, and Ali Farhadi. What can you learn from your muscles? learning visual representation from human interactions. *arXiv preprint arXiv:2010.08539*, 2020. [2](#)
- [14] Benjamin Elizalde, Soham Deshmukh, and Huaming Wang. Natural language supervision for general-purpose audio representations. In *ICASSP*. IEEE, 2024. [5, 6](#)
- [15] Jakob Engel, Kiran Somasundaram, Michael Goesele, Albert Sun, Alexander Gamino, Andrew Turner, Arjang Talatoff, Arnie Yuan, Bilal Souti, Brighid Meredith, et al. Project aria: A new tool for egocentric multi-modal ai research. *arXiv preprint arXiv:2308.13561*, 2023. [2, 4, 6, 7](#)
- [16] I-Sheng Fang and Jun-Cheng Chen. Camerabench: Benchmarking visual reasoning in mllms via photography. *arXiv preprint arXiv:2504.10090*, 2025. [2, 7, 3](#)
- [17] Daniel Geng, Charles Herrmann, Junhwa Hur, Forrester Cole, Serena Zhang, Tobias Pfaff, Tatiana Lopez-Guevara, Yusuf Aytar, Michael Rubinstein, Chen Sun, et al. Motion prompting: Controlling video generation with motion trajectories. In *CVPR*, 2025. [2](#)
- [18] Rohit Girdhar, Alaeldin El-Nouby, Zhuang Liu, Mannat Singh, Kalyan Vasudev Alwala, Armand Joulin, and Ishan Misra. Imagebind: One embedding space to bind them all. In *CVPR*, 2023. [1, 2](#)
- [19] Kristen Grauman, Andrew Westbury, Eugene Byrne, Zachary Chavis, Antonino Furnari, Rohit Girdhar, Jackson Hamburger, Hao Jiang, Miao Liu, Xingyu Liu, et al. Ego4D: Around the world in 3,000 hours of egocentric video. In *CVPR*, 2022. [4, 5](#)
- [20] Kristen Grauman, Andrew Westbury, Lorenzo Torresani, Kris Kitani, Jitendra Malik, Triantafyllos Afouras, Kumar Ashutosh, Vijay Baiyya, Siddhant Bansal, Bikram Boote, et al. Ego-Exo4D: Understanding skilled human activity from first-and third-person perspectives. In *CVPR*, 2024. [4, 6, 2](#)
- [21] Andrey Guzhov, Federico Raue, Jörn Hees, and Andreas Dengel. Audioclip: Extending clip to image, text and audio. In *ICASSP*. IEEE, 2022. [1, 2](#)
- [22] Richard Hartley and Andrew Zisserman. *Multiple view geometry in computer vision*. Cambridge university press, 2003. [1](#)
- [23] Masashi Hatano, Zhifan Zhu, Hideo Saito, and Dima Damen. The invisible egohand: 3d hand forecasting through egobody pose estimation. *arXiv preprint arXiv:2504.08654*, 2025. [2](#)
- [24] Hao He, Yinghao Xu, Yuwei Guo, Gordon Wetzstein, Bo Dai, Hongsheng Li, and Ceyuan Yang. Cameractrl: Enabling camera control for text-to-video generation. *arXiv preprint arXiv:2404.02101*, 2024. [2](#)
- [25] Jiahui Huang, Qunjie Zhou, Hesam Rabeti, Aleksandr Kordova, Huan Ling, Xuanchi Ren, Tianchang Shen, Jun Gao, Dmitry Slepichev, Chen-Hsuan Lin, et al. ViPE: Video pose engine for 3d geometric perception. *arXiv preprint arXiv:2508.10934*, 2025. [1, 2, 4, 6, 7, 5](#)

- [26] Hao Jiang and Kristen Grauman. Seeing invisible poses: Estimating 3d body pose from egocentric video. In *CVPR*. IEEE, 2017. 2
- [27] Nikhil Keetha, Norman Müller, Johannes Schönberger, Lorenzo Porzi, Yuchen Zhang, Tobias Fischer, Arno Knapitsch, Duncan Zauss, Ethan Weber, Nelson Antunes, et al. MapAnything: Universal feed-forward metric 3d reconstruction. *arXiv preprint arXiv:2509.13414*, 2025. 2
- [28] Leonid Keselman, John Iselin Woodfill, Anders Grunnet-Jepsen, and Achintya Bhowmik. Intel realsense stereoscopic depth cameras. In *CVPR workshops*, 2017. 2
- [29] Kris Kitani. Ego-action analysis for first-person sports videos. *IEEE Pervasive Computing*, 11(2), 2012. 2
- [30] Kris M Kitani, Takahiro Okabe, Yoichi Sato, and Akihiro Sugimoto. Fast unsupervised ego-action learning for first-person sports videos. In *CVPR*. IEEE, 2011. 2
- [31] Zhengfei Kuang, Shengqu Cai, Hao He, Yinghao Xu, Hongsheng Li, Leonidas J Guibas, and Gordon Wetzstein. Collaborative video diffusion: Consistent multi-video generation with camera control. *NeurIPS*, 37, 2024. 2
- [32] Gen Li, Yutong Chen, Yiqian Wu, Kaifeng Zhao, Marc Pollefeys, and Siyu Tang. EgoM2P: Egocentric multimodal multitask pretraining. *arXiv preprint arXiv:2506.07886*, 2025. 1, 2
- [33] Jiaman Li, Karen Liu, and Jiajun Wu. Ego-body pose estimation via ego-head pose estimation. In *CVPR*, 2023. 2
- [34] Ruilong Li, Brent Yi, Junchen Liu, Hang Gao, Yi Ma, and Angjoo Kanazawa. Cameras as relative positional encoding. *arXiv preprint arXiv:2507.10496*, 2025. 2
- [35] Xiaozhe Li, Kai Wu, Siyi Yang, YiZhan Qu, Guohua Zhang, Zhiyu Chen, Jiayao Li, Jiangchuan Mu, Xiaobin Hu, Wen Fang, et al. Can video generation replace cinematographers? research on the cinematic language of generated video. *arXiv preprint arXiv:2412.12223*, 2024. 2
- [36] Yin Li, Zhefan Ye, and James M Rehg. Delving into egocentric actions. In *CVPR*, 2015. 2
- [37] Zhengqi Li, Richard Tucker, Forrester Cole, Qianqian Wang, Linyi Jin, Vickie Ye, Angjoo Kanazawa, Aleksander Holynski, and Noah Snavely. MegaSaM: Accurate, fast and robust structure and motion from casual dynamic videos. In *CVPR*, 2025. 1, 2, 4, 6, 5
- [38] Kevin Qinghong Lin, Jinpeng Wang, Mattia Soldan, Michael Wray, Rui Yan, Eric Z Xu, Difei Gao, Rong-Cheng Tu, Wen-zhe Zhao, Weijie Kong, et al. Egocentric video-language pretraining. *NeurIPS*, 35, 2022. 1, 2, 5
- [39] Hongbo Liu, Jingwen He, Yi Jin, Dian Zheng, Yuhao Dong, Fan Zhang, Ziqi Huang, Yinan He, Yangguang Li, Weichao Chen, et al. ShotBench: Expert-level cinematic understanding in vision-language models. *arXiv preprint arXiv:2506.21356*, 2025. 2, 7, 3
- [40] Zhengyi Luo, Ryo Hachiuma, Ye Yuan, and Kris Kitani. Dynamics-regulated kinematic policy for egocentric pose estimation. *NeurIPS*, 34, 2021. 2
- [41] Niall Lyons, Avik Santra, and Ashutosh Pandey. Improved deep representation learning for human activity recognition using imu sensors. In *IEEE International Conference on Machine Learning and Applications (ICMLA)*. IEEE, 2021. 2
- [42] Lingni Ma, Yuting Ye, Fangzhou Hong, Vladimir Guzov, Yifeng Jiang, Rowan Postyeni, Luis Pesqueira, Alexander Gamino, Vijay Baiyya, Hyo Jin Kim, et al. Nymeria: A massive collection of multimodal egocentric daily motion in the wild. In *ECCV*. Springer, 2024. 4, 6
- [43] Seungwhan Moon, Andrea Madotto, Zhaojiang Lin, Aparajita Saraf, Amy Bearman, and Babak Damavandi. IMU2CLIP: language-grounded motion sensor translation with multimodal contrastive learning. In *Findings of the Association for Computational Linguistics: EMNLP*, 2023. 1, 2
- [44] Seungwhan Moon, Andrea Madotto, Zhaojiang Lin, Tushar Nagarajan, Matt Smith, Shashank Jain, Chun-Fu Yeh, Prakash Murugesan, Peyman Heidari, Yue Liu, et al. Anymal: An efficient and scalable any-modality augmented language model. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing: Industry Track*, 2024. 1, 2
- [45] Fangzhou Mu, Sicheng Mo, and Yin Li. Snag: Scalable and accurate video grounding. In *CVPR*, 2024. 1
- [46] Raul Mur-Artal, Jose Maria Martinez Montiel, and Juan D Tardos. Orb-slam: A versatile and accurate monocular slam system. *IEEE transactions on robotics*, 31(5), 2015. 2
- [47] Sanath Narayan, Mohan S Kankanhalli, and Kalpathi R Ramakrishnan. Action and interaction recognition in first-person videos. In *CVPR workshops*, 2014. 2
- [48] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018. 2, 3
- [49] Akhil Padmanabha, Saravanan Govindarajan, Hwanmun Kim, Sergio Ortiz, Rahul Rajan, Doruk Senkal, and Sneha Kadetotad. Egocharm: Resource-efficient hierarchical activity recognition using an egocentric imu sensor. *arXiv preprint arXiv:2504.17735*, 2025. 2
- [50] Hyun Soo Park, Jianbo Shi, et al. Force from motion: decoding physical sensation in a first person video. In *CVPR*, 2016. 2
- [51] Shraman Pramanick, Yale Song, Sayan Nag, Kevin Qinghong Lin, Hardik Shah, Mike Zheng Shou, Rama Chellappa, and Pengchuan Zhang. EgoVLPv2: Egocentric video-language pre-training with fusion in the backbone. In *ICCV*, 2023. 1, 2, 5, 6, 7, 3
- [52] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *ICML*. PmLR, 2021. 2, 3, 5, 6
- [53] Chris Rockwell, Joseph Tung, Tsung-Yi Lin, Ming-Yu Liu, David F Fouhey, and Chen-Hsuan Lin. Dynamic camera poses and where to find them. In *CVPR*, 2025. 4, 7, 2
- [54] Michael S Ryoo, Brandon Rothrock, and Larry Matthies. Pooled motion features for first-person videos. In *CVPR*, 2015. 2
- [55] Manolis Savva, Abhishek Kadian, Oleksandr Maksymets, Yili Zhao, Erik Wijmans, Bhavana Jain, Julian Straub, Jia Liu, Vladlen Koltun, Jitendra Malik, Devi Parikh, and Dhruv Batra. Habitat: A Platform for Embodied AI Research. In *ICCV*, 2019. 4

- [56] Johannes Lutz Schönberger and Jan-Michael Frahm. Structure-from-motion revisited. In *CVPR*, 2016. 2
- [57] Dian Shao, Yue Zhao, Bo Dai, and Dahua Lin. FineGym: A hierarchical video dataset for fine-grained action understanding. In *CVPR*, 2020. 4
- [58] Oriane Siméoni, Huy V Vo, Maximilian Seitzer, Federico Baldassarre, Maxime Oquab, Cijo Jose, Vasil Khalidov, Marc Szafraniec, Seungeun Yi, Michaël Ramamonjisoa, et al. Dinov3. *arXiv preprint arXiv:2508.10104*, 2025. 8
- [59] Suriya Singh, Chetan Arora, and CV Jawahar. Generic action recognition from egocentric videos. In *Fifth National Conference on Computer Vision, Pattern Recognition, Image Processing and Graphics (NCVPRIPG)*. IEEE, 2015. 2
- [60] Suriya Singh, Chetan Arora, and CV Jawahar. Trajectory aligned features for first person action recognition. *Pattern Recognition*, 62, 2017. 2
- [61] Khurram Soomro, Amir Roshan Zamir, and Mubarak Shah. UCF101: A dataset of 101 human actions classes from videos in the wild. *arXiv preprint arXiv:1212.0402*, 2012. 4
- [62] Maja Stikic, Tâm Huynh, Kristof Van Laerhoven, and Bernt Schiele. Adl recognition based on the combination of rfid and accelerometer sensing. In *international conference on pervasive computing technologies for healthcare*. IEEE, 2008. 2
- [63] Shuhan Tan, Tushar Nagarajan, and Kristen Grauman. Egodistill: Egocentric head motion distillation for efficient video understanding. *NeurIPS*, 36, 2023. 2
- [64] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *NeurIPS*, 30, 2017. 3
- [65] Heng Wang and Cordelia Schmid. Action recognition with improved trajectories. In *ICCV*, 2013. 2
- [66] Jianyuan Wang, Minghao Chen, Nikita Karaev, Andrea Vedaldi, Christian Rupprecht, and David Novotny. VGGT: Visual geometry grounded transformer. In *CVPR*, 2025. 1, 2, 4
- [67] Xinran Wang, Songyu Xu, Xiangxuan Shan, Yuxuan Zhang, Muxi Diao, Xueyan Duan, Yanhua Huang, Kongming Liang, and Zhanyu Ma. CineTechBench: A benchmark for cinematographic technique understanding and generation. *arXiv preprint arXiv:2505.15145*, 2025. 2
- [68] Yi Wang, Kunchang Li, Yizhuo Li, Yinan He, Bingkun Huang, Zhiyu Zhao, Hongjie Zhang, Jilan Xu, Yi Liu, Zun Wang, et al. InternVideo: General video foundation models via generative and discriminative learning. *arXiv preprint arXiv:2212.03191*, 2022. 1, 2
- [69] Yi Wang, Kunchang Li, Xinhao Li, Jiashuo Yu, Yinan He, Guo Chen, Baoqi Pei, Rongkun Zheng, Zun Wang, Yansong Shi, et al. InternVideo2: Scaling foundation models for multimodal video understanding. In *ECCV*. Springer, 2024. 1, 2
- [70] Yifan Wang, Jianjun Zhou, Haoyi Zhu, Wenzheng Chang, Yang Zhou, Zizun Li, Junyi Chen, Jiangmiao Pang, Chunhua Shen, and Tong He. Scalable permutation-equivariant visual geometry learning. *arXiv preprint arXiv:2507.13347*, 2025. 1, 2, 4, 6, 7, 5
- [71] Yuxi Xiao, Jianyuan Wang, Nan Xue, Nikita Karaev, Yuri Makarov, Bingyi Kang, Xing Zhu, Hujun Bao, Yujun Shen, and Xiaowei Zhou. Spatialtrackerv2: 3d point tracking made easy. *arXiv preprint arXiv:2507.12462*, 2025. 2
- [72] Fengyu Yang, Chao Feng, Ziyang Chen, Hyoungseob Park, Daniel Wang, Yiming Dou, Ziyao Zeng, Xien Chen, Rit Gangopadhyay, Andrew Owens, et al. Binding touch to everything: Learning unified multimodal tactile representations. In *CVPR*, 2024. 1, 2
- [73] Brent Yi, Vickie Ye, Maya Zheng, Yunqi Li, Lea Müller, Georgios Pavlakos, Yi Ma, Jitendra Malik, and Angjoo Kanazawa. Estimating body and hand motion in an ego-sensed world. In *CVPR*, 2025. 2
- [74] Ryo Yonetani, Kris M Kitani, and Yoichi Sato. Ego-surfing: Person localization in first-person videos using ego-motion signatures. *IEEE transactions on pattern analysis and machine intelligence*, 40(11), 2017. 2
- [75] Ye Yuan and Kris Kitani. 3d ego-pose estimation via imitation learning. In *ECCV*, 2018. 2
- [76] Mingfang Zhang, Yifei Huang, Ruicong Liu, and Yoichi Sato. Masked video and body-worn imu autoencoder for egocentric action recognition. In *ECCV*. Springer, 2024. 2
- [77] Zhengyou Zhang. Microsoft kinect sensor and its effect. *IEEE multimedia*, 19(2), 2012. 2
- [78] Yue Zhao, Ishan Misra, Philipp Krähenbühl, and Rohit Girdhar. Learning video representations from large language models. In *CVPR*, 2023. 1
- [79] Yi Zhou, Connelly Barnes, Jingwan Lu, Jimei Yang, and Hao Li. On the continuity of rotation representations in neural networks. In *CVPR*, 2019. 3, 5
- [80] Yang Zhou, Yifan Wang, Jianjun Zhou, Wenzheng Chang, Haoyu Guo, Zizun Li, Kaijing Ma, Xinyue Li, Yating Wang, Haoyi Zhu, et al. Omniworld: A multi-domain and multi-modal dataset for 4d world modeling. *arXiv preprint arXiv:2509.12201*, 2025. 6

# Seeing without Pixels: Perception from Camera Trajectories

## Supplementary Material

### 1. Supplementary Video

We invite readers to view the supplementary video available at <https://sites.google.com/view/seeing-without-pixels> for a visual demonstration of our work’s overview and additional qualitative examples. The video animates the static trajectories presented in the paper (*e.g.*, Fig. 1), offering a clearer view of the motion signatures. It also visualizes the learned embedding space, demonstrating how CamFormer clusters semantically similar neighbors. Furthermore, we provide qualitative examples of successful text retrieval (in both egocentric and exocentric domains) and demonstrate CamFormer’s emergent capabilities, such as repetitive action counting and text-based trajectory retrieval, alongside a visual analysis of failure modes.

### 2. Experimental Setup

#### 2.1. Task Setup

##### 2.1.1. Text Retrieval (5-way MCQ)

To evaluate cross-modal alignment between the camera trajectory and text modality, we formulate text retrieval as a 5-way MCQ task, mitigating the inherent ambiguity of open-ended retrieval. Given a trajectory query, the model must select the correct text description from five options based on feature similarity, using the pre-trained CamFormer as a frozen feature extractor.

**Negative Sampling Strategy.** To ensure a rigorous benchmark, we curate distractor options to have primary action verbs that do not overlap with the ground truth. Furthermore, we adopt a hard negative sampling strategy: distractor options are sourced from narrations within the same continuous video take (for egocentric datasets) or from captions with the same YouTube ID (for exocentric datasets), forcing the model to distinguish between temporally or thematically adjacent actions.

**Evaluation Splits.** This task serves as our primary testbed across Ego-Exo4D, Nymeria, and DynPose-100K. We further exploit dataset-specific annotations to dissect our model’s strengths versus the visual modality. On Ego-Exo4D, we analyze performance across “in-view” and “out-of-view” splits. On Nymeria, we break down results by narration type (legs, focus, body, hands). Qualitative MCQ examples are provided in Fig. 4 (main paper) and Fig. 8.

#### 2.1.2. Other Evaluation Tasks

**Proficiency Estimation on Ego-Exo4D.** We utilize the dataset’s skill-level annotations for the rock climbing and music activities. We focus exclusively on these two scenarios as they are the only ones that retain sufficient samples and a balanced class distribution after filtering for camera pose availability. To further mitigate imbalance, we formulate the task as a binary classification problem (expert vs. non-expert). We evaluate using an end-to-end fine-tuning protocol, comparing our pre-trained CamFormer initialization against a model trained from scratch.

**Keystep Recognition & Localization on Ego-Exo4D.**

We utilize the dataset’s established benchmark, which defines 278 fine-grained keystep labels. We adopt a frozen feature evaluation protocol for both tasks: we train a linear SVM for recognition and a dedicated localization network [45] for localization, both on top of our fixed trajectory embeddings.

**Activity Classification on Ego-Exo4D.** We evaluate coarse-grained understanding using the dataset’s 8 high-level activity labels. For this task, we adopt an end-to-end fine-tuning protocol, training the trajectory encoder jointly with a linear classification head. We compare initializing from our pre-trained CamFormer against a model trained from scratch to measure the benefit of pre-training.

**Scene Attribute Classification on DynPose-100K.** We enrich DynPose-100K with semantic labels to enable a new binary scene attribute classification task. To do this, we designed 10 questions covering a diverse range of attributes (*e.g.*, temporal, environmental, and social context) and utilized Gemini-2.5-Pro [7] to automatically label the videos. The prompts used for annotation are as follows:

1. **Day / Night:** Is the video filmed during the day or at night?
2. **Animal:** Does the video contain any animals?
3. **Text:** Does the video contain any visible written text?
4. **Urban / Rural:** Is the scene in the video urban or rural?
5. **Male / Female:** What is the gender of the people in the video?
6. **Food:** Does the video feature any food?
7. **Sports:** Is the video related to sports activities?
8. **Indoor / Outdoor:** Decide if the video is filmed indoors, outdoors, or if it is unclear.
9. **More than 1 people:** How many people are visible in the video?

## 10. Walking: Does the video show people walking?

From these new labels, we create a balanced 3,000-sample dataset for each attribute (with equal positive and negative examples) and train a linear SVM on our frozen camera trajectory features using an 80:20 train/test split.

**Event Classification on FineGym.** We evaluate on the 4 gymnasium event labels: floor exercise, balance beam, uneven bars, and vault. We adopt an end-to-end fine-tuning protocol, training the encoder jointly with a linear classification head.

**Action Recognition on UCF101-Dynamic.** To ensure a meaningful, motion-based evaluation, we curated this custom benchmark by quantitatively analyzing camera dynamics in UCF101. We selected the two most dynamic classes from each of four major action types, resulting in 8 classes: SkiJewandb sync ./anonymized/runs –project test-privacyt, SkateBoarding, Knitting, MoppingFloor, WalkingWithDog, Lunges, MilitaryParade, and SoccerPenalty. We evaluate using the same end-to-end fine-tuning protocol. For completeness, we also report results on the full 101-class UCF101 dataset (cf. Table 8).

## 2.2. Datasets

**Pretraining Data.** For the egocentric domain, we use Ego-Exo4D [20]. Adhering to the official splits, we obtain 159,186 training and 69,073 validation (trajectory, text) pairs, where text is human-annotated narrations provided by the dataset. Trajectory boundaries for these long, untrimmed videos are defined following [38, 51]. There are two camera trajectory sources: the original Aria glasses data (which we downsample to 20 FPS) and video-only estimations we obtain by running  $\pi^3$  (5 FPS); a comparison is provided in Table 10. For the exocentric domain, we use DynPose-100K [53]. As no official split is available, we randomly split the dataset into 88,151 training and 10,452 validation (trajectory, text) pairs, where the text is video captions from Panda70M [6]. Unlike the egocentric data, these are short clips with fixed boundaries. There are two trajectory sources, both estimated from videos: the original dataset’s provided ones (12 FPS) and the ViPE-provided [25] ones (30 FPS).

**Downstream Data.** For the egocentric domain, we evaluate on Ego-Exo4D and Nymeria. On Ego-Exo4D, for our designed text retrieval task, we draw samples from the official validation split; for all other tasks, we follow the dataset’s official benchmark splits. For Nymeria, which shares the same Aria hardware as Ego-Exo4D, we use the entire set of data with available motion narrations as a zero-shot test set for text retrieval; the camera trajectory is also

downsampled to 20FPS. For the exocentric domain, we use the original action labels from FineGym and UCF101. Since these datasets lack trajectories, we generate them using our pose estimation pipeline at 5 FPS. We generate three versions for UCF101 (using MegaSam [37], ViPE [25], and  $\pi^3$  [70]) and one version for FineGym (using  $\pi^3$ ).

## 2.3. Implementation

When using hardware-estimated camera poses (*i.e.*, from Aria glasses [15]), we utilize the gravity direction information. To be specific, we compute the 3D gravity vector in our chosen relative reference frame and project it to  $d_{in}$  via a learned linear projection layer. This additional token is subsequently prepended to the input sequence before processing by the Transformer. This step is omitted when processing poses estimated from monocular videos. We train CamFormer using an AdamW optimizer with a learning rate of  $1 \times 10^{-4}$ , weight decay of  $1 \times 10^{-3}$ , and a batch size of 1024. CLIP loss temperature  $\tau$  is 0.07. Training is conducted on 8 NVIDIA A100 GPUs.

**Camera Pose Extraction.** Given the computational expense of running multiple pose estimators, all pose estimations are performed at 5 FPS. For the particularly long, untrimmed videos in the Ego-Exo4D activity classification task (which can be several minutes), we further limit pose extraction to the center 4-second clip. For the egocentric setting, an alignment step is required. Our model is pre-trained on the Aria pose coordinate frame (x left, y up, z forward), so we apply a rigid transformation to convert all estimated poses from the standard OpenCV frame (x right, y down, z forward) before they are fed into our model.

**Analysis of Contextualized Trajectory Encoding.** We detail the experimental setup for the four settings in Fig. 6 of the main paper.

- *Global Label Tasks.* We fine-tune CamFormer end-to-end for these tasks. Our model accepts flexible input trajectory lengths (both in pre-training and finetuning), which allows us to systematically vary the input length during inference. (1) For Ego-Exo4D activity classification, we vary the input trajectory length from 1 to 16 seconds. (2) For FineGym event classification, since events have fixed segments, we vary the input ratio from 20% to 100% of the full event duration.
- *Localized Label Tasks.* For these tasks, we apply the pre-trained CamFormer as a frozen feature extractor and investigate extending the context outside the given segment window  $[t_1, t_2]$ . (3) For Ego-Exo4D text retrieval, we extend the atomic window  $[t_1, t_2]$  by a total duration of 0, 2, 4, 6, or 8 seconds. We apply this symmetrically; for instance, a 2-second extension results in the input window  $[t_1 - 1, t_2 + 1]$ . (4) For Ego-Exo4D keystep recogni-

Table 5. Per-Activity Breakdown of Ego-Exo4D Text Retrieval Results (MCQ Accuracy). We compare the performance of CamFormer (trajectory features) and the best video baseline, EgoVLPv2 [51] (Ego-Exo4D), across all 8 activity scenarios. This detailed breakdown allows us to identify the relative strengths and weaknesses of the camera trajectory versus the video modality.

	Bike		Health		Cooking		Music		Soccer		Dance		Basketball		Bouldering		All
	iv	oov															
# Queries	500	205	500	416	500	500	500	176	500	500	282	500	500	500	500	500	7079
Video	<b>37.20</b>	<b>40.98</b>	<b>52.60</b>	<b>47.12</b>	<b>61.60</b>	<b>45.80</b>	29.40	21.02	24.40	21.20	25.89	15.40	50.80	38.00	59.40	29.60	38.40
Ours	28.80	28.78	24.40	26.44	49.80	39.40	<b>32.00</b>	<b>35.23</b>	<b>52.00</b>	<b>26.60</b>	<b>39.36</b>	<b>50.60</b>	<b>67.40</b>	<b>57.00</b>	<b>62.60</b>	<b>55.40</b>	<b>44.80</b>

tion, we expand the input trajectory window proportionally (100%-400% of the original duration).

## 2.4. Baselines

(1) Direct LMM Baseline. For the Gemini row in Table 2 of the main paper, we input 8 uniformly sampled video frames directly. The five candidate texts are randomly shuffled and assigned labels (A-E). Gemini-2.5-Pro [7] is queried with the prompt: *Which of the following descriptions best matches the video?*

(2) Camera Description Baseline. For exocentric text retrieval (Table 4 of the main paper), we adopt a two-stage process. First, we prompt specialized LMMs (Qwen-VL-7B [16] or ShotVL-7B [39]) to *Describe the camera motion in this video*. Note that these models are trained to generate textual description of the camera motion in the associated video (e.g., “zoom”, “pan”) and do not aim to describe the content of the video. Second, we feed this generated description into another LLM (we use Gemini-2.5-Flash [7]) and prompt it to answer the MCQ given the camera motion described in text form. The prompt is: *The following describes the motion and focus of a camera while filming a scene:[Generated Description]. Which of the following events or scene descriptions is most likely being filmed with this camera movement? [Option A-E]*.

## 3. Results

### 3.1. Additional Results

**Text Retrieval.** Supplementing Table 2 of the main paper, Table 5 provides a detailed activity-level breakdown on Ego-Exo4D text retrieval, comparing our CamFormer embeddings with leading video encoder features (EgoVLPv2 [51]). The results allow us to clearly delineate the strengths of the two modalities. For the procedural activities (where visual cues are more critical), the video baseline maintains its lead. For the five physical activities, the camera trajectory modality is demonstrably stronger on its own. This performance is particularly effective in out-of-view (oov) settings, highlighting trajectory’s unique value in scenarios where the visual signal is occluded or ambiguous. Finally, our CamFormer achieves the best overall result, demonstrating the promising value of the camera tra-

Table 6. Proficiency Estimation Accuracy (%) on Ego-Exo4D. The camera trajectory is a particularly strong modality for this physical task, outperforming the video baseline. Moreover, our pretraining is crucial, as initializing from CamFormer provides a boost over training from scratch.

Method	Modality	Pretrain?	Bouldering	Dancing
Majority	-	-	55.97	59.30
TimeSformer [4]	video	-	55.35	69.92
CamFormer	trajectory	✗	63.52	66.67
CamFormer	trajectory	✓	<b>65.41</b>	<b>70.73</b>

Table 7. Keystep Recognition and Localization Results on Ego-Exo4D. For these procedural tasks, where vision is a strong baseline, fusing trajectory and video features (denoted by \*) consistently outperforms the video-only model, proving that camera trajectory provides an essential, non-redundant signal.

Method	Modality	Rec.	Loc.	Rank@1	Loc.	Rank@5
		Acc.	IoU	IoU	IoU	IoU
Majority	-	3.52	-	-	-	-
EgoVLPv2 [51]	video	29.17	31.81	26.28	62.90	52.69
CamFormer	trajectory	14.07	20.29	15.67	47.23	38.09
CamFormer*	video+trajectory	<b>32.37</b>	<b>34.68</b>	<b>29.06</b>	<b>66.65</b>	<b>57.29</b>

jectory in action understanding.

**Qualitative Results.** Supplementing Fig. 4 of the main paper, Fig. 8 presents additional qualitative results for the text retrieval task. These examples demonstrate that across diverse datasets and varying text descriptions, CamFormer successfully decodes the semantic information embedded in camera trajectories and accurately matches it with the corresponding text.

**Proficiency Estimation.** For physical activities, we find that camera trajectory is a powerful standalone signal for assessing skill levels (beginner/expert). Table 6 presents proficiency estimation results for the two physical activities. Our lightweight model successfully captures the motion signatures of expertise and outperforms the video baseline. This performance is further boosted by our pre-training strategy, which surpasses the train-from-scratch counterpart.

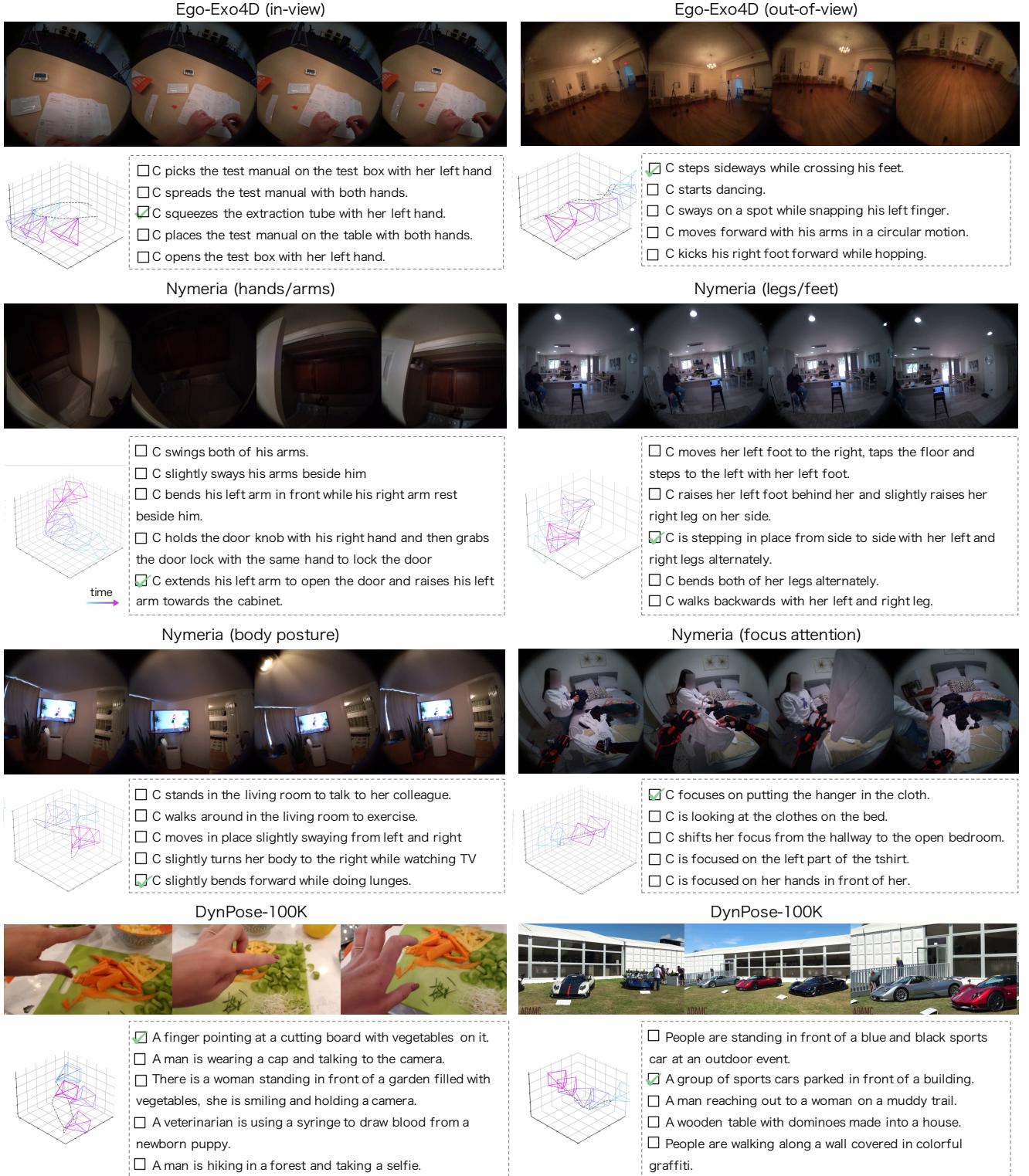


Figure 8. Qualitative Text Retrieval Results on Ego-Exo4D (top row), Nymeria (middle two rows) and DynPose-100K (bottom row). Note that CamFormer takes only the camera trajectory as input; corresponding video frames are shown solely for illustration. These examples further demonstrate that our model effectively captures the trajectory-semantic link across both egocentric and exocentric domains.

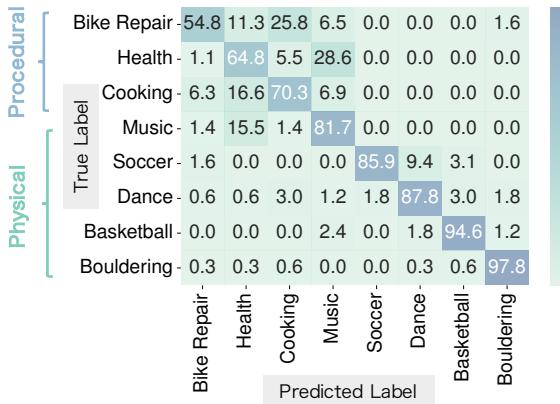


Figure 9. Confusion Matrix for CamFormer on Ego-Exo4D Activity Classification. CamFormer performs strongly on dynamic physical activities (*e.g.*, near-perfect for ‘‘bouldering’’), while the main confusion occurs between the three procedural activities, which involve more subtle motion cues.

Table 8. Comparing action recognition results on UCF101-Dynamic (left) and UCF101 (right) with various estimated camera poses. Echoing our egocentric analysis, the results confirm the benefits of our pre-training strategy. Across all pose estimators, the model initialized with our checkpoint pre-trained on DynPose-100K (✓) consistently outperforms its counterpart trained from scratch (✗).

Pose Source	UCF101-Dynamic			UCF101		
	Pretrain ✗	Pretrain ✓	Δ	Pretrain ✗	Pretrain ✓	Δ
	66.67	69.15	+2.48	16.02	17.91	+1.89
MegaSaM [37]	64.18	68.16	+3.98	16.54	19.62	+3.08
ViPE [25]	61.69	64.18	+2.49	17.53	19.25	+1.72
π <sup>3</sup> [70]						

**Keystep Recognition & Localization.** For procedural activities, where motion patterns can be ambiguous, camera trajectory provides valuable complementary information. For the keystep recognition and localization tasks on these scenarios (Table 7), fusing trajectory with video features provides a consistent performance boost over the video-only baseline.

**Activity & Event Classification.** The confusion matrix (Fig. 9) on Ego-Exo4D activity classification provides a detailed breakdown of the per-class accuracy results in Fig. 3 (c). CamFormer excels at recognizing dynamic physical activities, while confusion is heavily concentrated among the three procedural activities (where camera motion is subtle). For the exocentric domain, on FineGym event classification (Fig. 10), the model performs strongly on recognizing ‘‘floor exercise’’, and the main confusion occurs between ‘‘uneven bars’’ and ‘‘balance beam’’.

**Action Recognition.** Table 8 compares action recognition performance using our pre-trained CamFormer against

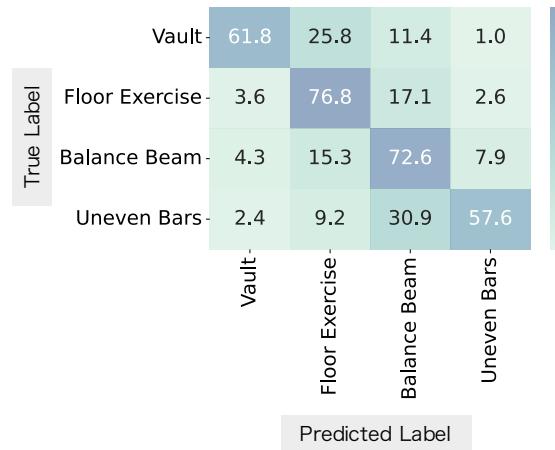


Figure 10. Confusion Matrix for CamFormer on FineGym Event Classification. The matrix details our model’s performance on the 4 gymnasium activities.

Table 9. Ablation Study of Input Camera Trajectory Representation on Ego-Exo4D Text Retrieval. We compare various formulations, including the use of absolute vs. relative poses, the rotation format (no vs. 4D quaternion vs. 6D continuous), the specific reference frame used for calculating relative poses, and whether to include gravity direction information.

	Dim. (Tsl. + Rot.)	Acc. (%)
Absolute	3D	32.84
Absolute	3D + 4D	34.74
Absolute	3D + 6D	37.82
Relative (prev.)	3D + 4D	43.66
Relative (mid.)	3D + 4D	44.02
Relative (any)	3D + 4D	44.00
Relative (mid.)	3D + 6D	44.12
+ Gravity direction	3D + 6D	<b>44.81</b>

the train-from-scratch baseline on UCF-Dynamic (our curated 8-class subset) and the full UCF101 dataset. The results show that initializing from CamFormer yields better performance than the train-from-scratch baseline across all three pose sources and on both settings. The largest performance gain is observed when using ViPE poses, as CamFormer was pre-trained with ViPE camera trajectories on DynPose-100K. Even with the other pose sources, the consistent gains observed across datasets demonstrate the robust generalization capability of our pre-training strategy.

### 3.2. Ablation Study

Table 9 presents the ablation study on Ego-Exo4D text retrieval, where we compare various ways to represent the input camera trajectory. The results demonstrate that relative pose sequences are critical and greatly outperform absolute pose sequences, with the sequence midpoint being the optimal reference frame. Furthermore, the 6D continuous rotation representation [79] is preferred over the 4D quater-

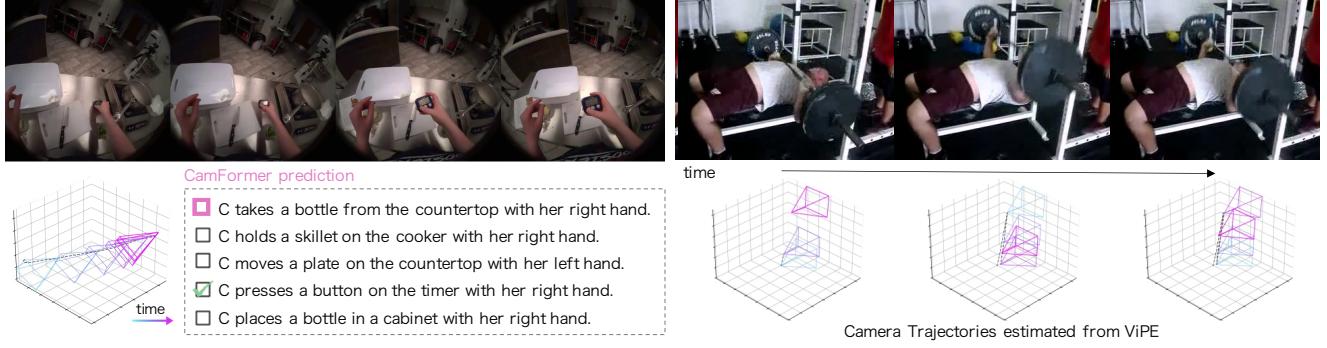


Figure 11. Failure Cases. Left: CamFormer struggles to distinguish actions with subtle motion patterns and often cannot capture specific noun semantics. Right: The pose source fails when the ViPE estimator mistakenly tracks object motion as camera motion (best viewed in Supp. video), highlighting a necessary area for development in pose estimation algorithms.

Table 10. Ablation Study of Pretraining Choices on Ego-Exo4D text Retrieval. We compare CamFormer performance using two different camera pose sources (high-fidelity Aria [15] vs. video-estimated  $\pi^3$  [70]) and text encoder training modes (frozen vs. finetuned).

	Pose Source	Text Encoder	Acc. (%)
(a)	$\pi^3$ [70]	frozen	43.86
(b)	Aria [15]	finetune	45.42
(c)	Aria [15]	frozen	44.81

nion, and encoding gravity direction provides a further performance boost.

Table 10 investigates our pretraining choices. First, regarding pose source: replacing high-fidelity Aria poses with video-estimated  $\pi^3$  ones still yields comparable performance (43.86% vs. 44.81%). This is a promising result, indicating significant potential to scale up pre-training data using poses estimated from large collections of in-the-wild videos. Second, regarding the text encoder: while fine-tuning the CLIP encoder yields a marginal performance gain, it comes with a substantial computational cost. We therefore adopt the frozen text encoder for our final model to prioritize efficiency, though we posit that end-to-end fine-tuning may become more beneficial as data scale increases in the future.

### 3.3. Limitations

We acknowledge that obtaining high-quality camera poses initially incurs a computational cost, whether through multi-sensor hardware or video estimation algorithms. We view this, however, as a one-time, amortized process. Concurrent advances in hardware and the development of efficient algorithms are actively enriching existing video datasets with camera trajectories. This growing repository of pose-annotated data, like [80], provides the reusable, large-scale foundation that our method can directly leverage.

Due to the inherent differences between egocentric and

exocentric motion, we currently train a separate CamFormer for each domain under our unified framework. A promising avenue for future work is to build a single, unified trajectory encoder. This could be achieved by introducing an explicit conditional domain token that allows the unified encoder to effectively distinguish and interpret the recorder’s intent across both camera perspectives.

Finally, we present two failure modes of our investigation, as shown in Fig. 11. The left panel reveals an intrinsic limitation of the camera trajectory signal: it struggles with subtle motion patterns (“press a button” in this case), confuses it with the adjacent action of “taking something from the countertop”, and inherently fails to encode specific noun semantics. The right panel highlights an issue with the pose source: we observe cases where the estimator (*e.g.*, ViPE [25]) mistakenly correlates object motion with camera motion. This failure suggests that further algorithmic development in camera pose estimation is necessary to ensure robust semantic analysis.