

# EgoPrompt: Prompt Learning for Egocentric Action Recognition

Huaihai Lyu

MAIS, Institute of Automation,  
Chinese Academy of Sciences  
Beijing, China  
lvhuaihai2023@ia.ac.cn

Yuheng Ji

MAIS, Institute of Automation,  
Chinese Academy of Sciences  
Beijing, China  
jiyuheng2023@ia.ac.cn

Chaofan Chen\*

MAIS, Institute of Automation,  
Chinese Academy of Sciences  
Beijing, China  
chencfbupt@gmail.com

Changsheng Xu

MAIS, Institute of Automation,  
Chinese Academy of Sciences and Peng Cheng Laboratory  
Beijing, China  
csxu@nlpr.ia.ac.cn

## Abstract

Driven by the increasing demand for applications in augmented and virtual reality, egocentric action recognition has emerged as a prominent research area. It is typically divided into two subtasks: recognizing the performed behavior (i.e., *verb component*) and identifying the objects being acted upon (i.e., *noun component*) from the first-person perspective. However, most existing approaches treat these two components as independent classification tasks, focusing on extracting component-specific knowledge while overlooking their inherent semantic and contextual relationships, leading to fragmented representations and sub-optimal generalization capability. To address these challenges, we propose a prompt learning-based framework, **EgoPrompt**, to conduct the egocentric action recognition task. Building on the existing prompting strategy to capture the component-specific knowledge, we construct a **Unified Prompt Pool** space to establish interaction between the two types of component representations. Specifically, the component representations (from verbs and nouns) are first decomposed into fine-grained patterns with the prompt pair form. Then, these pattern-level representations are fused through an attention-based mechanism to facilitate cross-component interaction. To ensure the prompt pool is informative, we further introduce a novel training objective, **Diverse Pool Criteria**. This objective realizes our goals from two perspectives: *Prompt Selection Frequency Regularization* and *Prompt Knowledge Orthogonalization*. Extensive experiments are conducted on the Ego4D, EPIC-Kitchens, and EGTEA datasets. The results consistently show that EgoPrompt achieves state-of-the-art performance across within-dataset, cross-dataset, and base-to-novel generalization benchmarks.

## CCS Concepts

• **Computing methodologies** → *Scene understanding*; **Cross-validation**; **Learning paradigms**.

\*Corresponding author.



This work is licensed under a Creative Commons Attribution 4.0 International License.  
MM '25, Dublin, Ireland

© 2025 Copyright held by the owner/author(s).  
ACM ISBN 979-8-4007-2035-2/2025/10  
<https://doi.org/10.1145/3746027.3754749>

## Keywords

Egocentric action recognition; domain generalization; prompt tuning

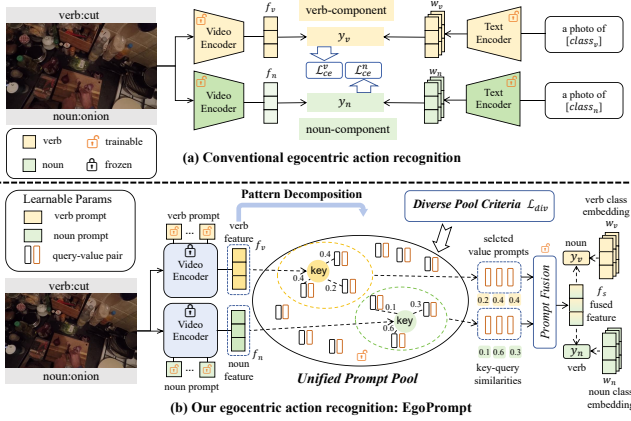
## ACM Reference Format:

Huaihai Lyu, Chaofan Chen, Yuheng Ji, and Changsheng Xu. 2025. EgoPrompt: Prompt Learning for Egocentric Action Recognition. In *Proceedings of the 33rd ACM International Conference on Multimedia (MM '25)*, October 27–31, 2025, Dublin, Ireland. ACM, New York, NY, USA, 9 pages. <https://doi.org/10.1145/3746027.3754749>

## 1 Introduction

Recent advancements in augmented and virtual reality (AR/VR) technologies [1, 2] have demonstrated the potential to transform the way humans interact with the digital world. A fundamental requirement for such AR/VR systems is the ability to recognize user behaviors (*verb component*) and the objects they interact with (*noun component*) from egocentric videos. This demand has spurred a growing interest in egocentric action recognition (EAR) [3–5], which focuses on understanding first-person visual data to enable applications such as assistive systems and wearable devices.

Although existing EAR methods [6, 7] have made significant progress, their generalization performance remains insufficient when adapting to real-world datasets with distribution shifts. Such distribution shifts arise from variations in the environment, recording conditions, and user behaviors, posing significant challenges for robust egocentric video understanding [8–12]. To address this issue, researchers have introduced some parameter-efficient strategies [13], such as prompt learning [14–16] and adapter-based techniques [17, 18], to effectively adapt to the EAR tasks. Unlike third-person perspectives that require a comprehensive understanding of the entire environmental context, the semantic content of egocentric video is primarily reflected in the Human Object Interaction (HOI) region [6]. However, the visual field in egocentric videos is often cluttered with *irrelevant objects and background noise* in the HOI region, which significantly complicates the action recognition task. As shown in Fig. 1 (a), most existing methods treat verbs and nouns as separate entities to capture component-specific knowledge. These methods neglect the strong relationship between behaviors and the objects being acted upon in EAR. For example, the interactive object “carrot” could constrain the semantic space of verbs (e.g., “slicing” and “cutting”) depending on the context, and the



**Figure 1: Comparison with existing framework. (a) Conventional methods use the component label of verb/noun to fine-tune its corresponding encoder independently. (b) Our proposed EgoPrompt constructs a *Unified Prompt Pool* with *Diverse Pool Criteria* constraint, which decomposes component-specific representation into implicit prompt pair patterns and achieves better knowledge interaction with an attention fusion mechanism.**

verb “cutting” has specified the object of interaction with attributions (e.g., “cuttable” and “solid”). This interdependence highlights that they are not independent components but instead mutually influence one another, which collectively reflects the egocentric semantic content and helps in identifying HOI information from the noisy environment.

Based on the above analysis, we propose a novel prompt learning-based framework, **EgoPrompt**, to explore the nature of semantic interplay in EAR task. Built upon the component-specific feature extraction, we further construct a **Unified Prompt Pool** space that encodes fine-grained implicit patterns through multiple query-value prompt pairs. Specifically, as illustrated in Fig 1 (b), the component-specific representations serve as keys that match with queries in the prompt pool. Matched queries retrieve corresponding value prompts, representing the decomposed latent patterns from the input representation. Finally, the selected value prompts are integrated into a fused representation that facilitates cross-component semantic interaction. By enabling the verb and noun features to interact with pattern prompt pairs through attention-based fusion, the egocentric model effectively captures their contextual semantic interplay. This interaction facilitates the integration of complementary information from both components, leading to a more context-aware understanding of HOI semantics. To ensure the Unified Prompt Pool captures the informative pattern knowledge, we introduce a novel training objective, **Diverse Pool Criteria**, to consider the following two factors: (1) *Prompt Selection Frequency Regularization*, ensuring balanced utilization of prompts by discouraging overused ones and encouraging seldom-selected prompts, and (2) *Prompt Knowledge Orthogonalization*, minimizing redundancy and fostering semantic diversity through reducing the cosine similarities between prompt pairs. This mechanism can help the prompt pool capture a wide range of contextual cues, thereby improving the quality of the learned video representations.

To sum up, our main contributions are fourfold:

- We propose a novel prompt learning framework, **EgoPrompt**, tailored for egocentric action recognition. It explores the effect of semantic interplay in HOI understanding, addressing the unique challenges in egocentric video understanding.
- We introduce the **Unified Prompt Pool** design, enabling cross-component interaction and capturing the HOI semantic information.
- We present a training objective, **Diverse Pool Criteria**, to enhance prompt informativeness by encouraging balanced usage and enforcing semantic orthogonality.
- We conduct extensive experiments, including **within/cross-dataset and base-to-novel generalization**. Results consistently demonstrate the effectiveness of EgoPrompt in improving generalization capability.

## 2 Related Work

### 2.1 Prompt Learning

Prompt learning [17, 19–21] has emerged as a lightweight and parameter-efficient alternative to traditional fine-tuning, enabling task-specific adaptations without modifying the entire backbone. Initially, prompt learning techniques were applied in the textual modality, where handcrafted prompts such as “a photo of a [CLASS]” were embedded in models like CLIP [22] for zero-shot classification tasks. However, hand-crafted prompts often lack adaptability and explainability. To address these limitations, CoOp [19] introduced continuous prompts represented as trainable vectors appended to text tokens, while CoCoOp [20] further proposed image-conditioned prompts at the instance level to improve novel class generalization [23]. To align the representations of dual-encoder architectures, PromptSRC [24] jointly tunes visual and textual encoders, employing regularization constraints between original and adapted representations. MaPL [16] extends prompt learning to dual encoders by introducing deep prompts and a text-to-visual coupling module, effectively projecting textual knowledge into the visual encoder. Similarly, KgCoOp [25] and TCP [15] emphasize class-aware knowledge encoded in the text encoder to enhance generalization. L2P [26] models the task-specific knowledge with a prompt-pair form, providing a solution to catastrophic forgetting. Due to its lightweight advantage, prompt learning has started to gain attention in EAR. For example, POV [7] employs hand-object interaction (HOI) labels to guide optimization and uses exocentric datasets for auxiliary learning. However, the application of prompt learning to egocentric action recognition remains underexplored. In this work, we introduce prompt learning as a way to enhance the generalization capability of the egocentric model. By integrating techniques such as deep prompts and prompt pools, our proposed EgoPrompt framework enhances generalization capabilities while addressing the unique characteristics of egocentric video data.

### 2.2 Egocentric Action Recognition

The emergence of large-scale egocentric datasets such as Ego4D [10], EGTEA [9], and Epic-Kitchens [11] has significantly advanced research in egocentric video understanding [27, 28]. These datasets provide diverse benchmarks for exploring egocentric perception, including action recognition [6, 29], video summarization [3, 4, 30],

and object interaction understanding [31, 32]. To address the practical demands of egocentric action recognition, researchers have proposed various methods to enhance the efficiency and generalization of egocentric models. For instance, EgoDistill [3] introduces IMU signals to reduce the computational cost of processing dense keyframes, while Ego-Only [33] explores differences between egocentric and exocentric videos, employing a masked autoencoder [34] during pretraining to improve video understanding capabilities. Despite these advancements, most existing methods focus on intra-dataset evaluation, which limits their applicability to real-world scenarios where distribution shifts are prevalent. X-MIC [29] extends the generalization scenarios proposed in [18, 35, 36] for egocentric action recognition, including cross-dataset and base-to-novel generalization. Additionally, AoP [6] decouples egocentric action recognition into component-aware tasks, introducing an adapter-based approach that uses verb embeddings as prior knowledge to assist noun prediction under an open-vocabulary setting [37]. While these approaches demonstrate the importance of component-aware modeling, they often neglect the constraints that nouns impose on verbs, which may hinder generalization performance. Nonetheless, the potential of fully exploring component-based representation fusion remains untapped. Current methods fail to capture the semantic interplay between verbs and nouns, which is crucial for the goal-oriented nature (focusing on the HOI region) of egocentric actions. To address this challenge, we propose EgoPrompt, which incorporates a Unified Prompt Pool design to enable interaction between verb and noun component-specific representations.

### 3 Methodology

In this section, we introduce the technical details of our proposed **EgoPrompt**, a novel prompt learning approach specifically designed for the EAR task. Building on the foundations of the EAR baseline (described in Sec. 3.1), EgoPrompt further facilitates the component knowledge interaction (described in Sec. 3.2), achieving a better understanding of egocentric semantic information.

#### 3.1 Baseline

In this work, we follow X-MIC [29] to adopt LaVILA [8] as the backbone for egocentric action recognition. LaVILA is a dual-encoder vision-language model that processes video and text inputs through separate transformer-based encoders. To enhance cross-modal alignment, we incorporate the deep prompt learning strategy from MaPLe [16]. In this design, learnable textual prompts are inserted into each transformer layer of the text encoder. These prompts are then transformed into video-aware prompts via lightweight layer-wise projection modules and injected into the corresponding layers of the video encoder. For implementation details, please refer to MaPLe [16] and our supplementary material.

To adapt to the EAR task, we feed two sets of video prompts into the video encoders to extract component-specific representations  $f_v$  and  $f_n$  for verbs and nouns. In the same way, the text encoder processes component-specific text prompts to generate the corresponding class embeddings  $W^v, W^n$ . These text prompts are initialized using the tokenized results of our selected hand-crafted templates. We empirically study the impact of template choice

in the supplementary material and adopt the most effective ones, as illustrated in the Template Initialization stage of Fig. 2. After this, both visual and textual prompts are jointly optimized using a component-specific classification loss:

$$\mathcal{L}_{ce}^c = \frac{\exp(\cos(f_c, w_y^c)/\tau)}{\sum_{j=1}^{N^c} \exp(\cos(f_c, w_j^c)/\tau)}, \quad c \in \{v, n\}, \quad (1)$$

where  $c$  denotes the component between verb and noun, and  $w_y^c$  is the text embedding of the ground-truth class  $y$  for component  $c$ .  $\cos(\cdot)$  denotes cosine similarity,  $\tau$  is a temperature factor, and  $N^c$  is the number of classes for component  $c$ .

In addition, we apply a knowledge-guided loss [25] to preserve the semantic priors encoded in hand-crafted templates:

$$\mathcal{L}_{kg}^c = \|W^c - \hat{W}^c\|_2^2, \quad c \in \{v, n\}, \quad (2)$$

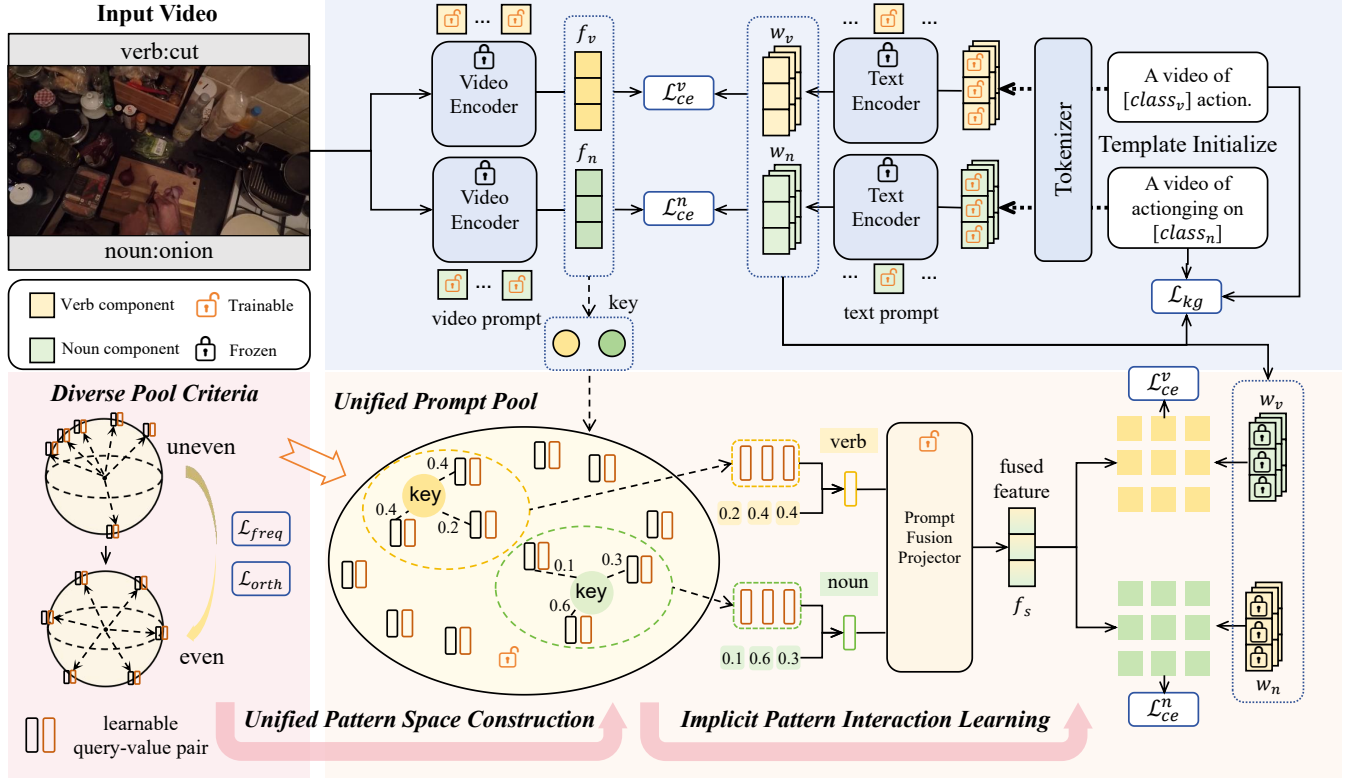
where  $\hat{W}^c$  denotes the class embedding obtained directly from the tokenized hand-crafted templates (without further training), while  $W^c$  is derived from text prompts that are initialized from these templates and then optimized during training. This loss encourages the learned soft prompts to be semantically consistent with the original template.

#### 3.2 EgoPrompt

As described above, the baseline lacks the ability to reason over the inherent correlations between actions and objects, which is essential for understanding egocentric HOI knowledge. Building upon the component-specific knowledge, we further develop **EgoPrompt** to facilitate the interaction within components. Specifically, EgoPrompt introduces two key modules: the *Implicit Pattern Interaction Learning* and *Unified Pattern Space Construction* strategy.

**3.2.1 Implicit Pattern Interaction Learning.** Although the component-specific prompting design effectively encodes component semantics, egocentric models still face challenges in irrelevant background noise and object clutter in egocentric scenes. Fortunately, the strong contextual association between verbs and nouns offers valuable cues for disambiguating such noisy scenes and identifying the core HOI knowledge. To fully exploit this semantic interplay, we introduce a **Unified Prompt Pool**, which serves as a shared latent space for capturing implicit cross-component patterns. As illustrated in Fig. 2, this module enables dynamic interaction between the component-specific features by decomposing them into fine-grained patterns and re-integrating them into a unified representation.

The interaction process is achieved by utilizing the component-specific features  $f_v$  (verb feature) and  $f_n$  (noun feature) as keys to interact with query-value pairs in the Unified Prompt Pool. Specifically, the prompt pool contains  $P$  prompt pairs, where each pair consists of a query prompt  $q$  and a value prompt  $v$ , denoted as:  $\mathcal{P} = \{(q^1, v^1), (q^2, v^2), \dots, (q^P, v^P)\}$ . The query prompt serves as the connection between the component-specific representation space and the unified prompt pool space, and the value prompt encodes the corresponding latent semantic pattern. For each component (verb and noun), we first retrieve the Top- $k$  most semantically similar query prompts based on the cosine similarity metric. Then, we compute the attention weights  $\alpha_i^c$  over the selected prompts



**Figure 2: Overall framework of EgoPrompt.** Building upon the Baseline work, EgoPrompt further establishes the semantic interaction between components. Specifically, under the guidance of the Diverse Pool Criteria, EgoPrompt constructs the Unified Pattern Space with the Unified Prompt Pool design. It decomposes the component-specific representation into fine-grained patterns and integrates the selected prompt pairs into a fused representation.

using a temperature-scaled softmax:

$$\alpha_i^c = \frac{\exp(\cos(f_c, q_i)/\tau)}{\sum_{q_j \in Q_c} \exp(\cos(f_c, q_j)/\tau)}, \quad q_i \in Q_c = \text{Top-}k(f_c, Q), \quad (3)$$

where  $Q$  denotes the query prompt set. Accordingly, we fetch the corresponding value prompts  $\mathcal{V}_c = \{v_i | q_i \in Q_c\}$  and compute the pattern-composed feature for each component as:

$$f'_c = \sum_{v_i \in \mathcal{V}_c} \alpha_i^c \cdot v_i, \quad c \in \{v, n\}. \quad (4)$$

After this, the final fusion representation is obtained by projecting the two pattern-based component representations:

$$f_s = \text{Proj}(f'_v, f'_n), \quad (5)$$

where **Proj**( $\cdot$ ) is a learnable projection layer. This fusion representation  $f_s$  encodes the implicit interaction between verbs and nouns, providing a noise-robust and component-shared embedding well aligned with egocentric HOI recognition.

**3.2.2 Unified Pattern Space Construction.** To ensure that the prompt pool encodes a rich variety of informative patterns, we introduce the **Diverse Pool Criteria**, a dedicated objective function designed to promote both effective prompt utilization and semantic diversity

among prompt pairs. This objective comprises two complementary regularization strategies:

1) **Prompt Selection Frequency Regularization:** To avoid prompt over-reliance and encourage balanced utilization across the pool, we reward the  $k$  least-frequently selected prompts (encouraging more utilization) and penalize the  $k$  most-frequently selected prompts (discouraging over-reliance). Let  $c_p^c$  denote the selection count of prompt pair  $p$  for component  $c \in \{v, n\}$  during training. The regularization term is defined as:

$$\mathcal{L}_{\text{freq}} = - \sum_{c \in \{v, n\}} \left( \sum_{p \in S_{\min}} c_p^c - \sum_{p \in S_{\max}} c_p^c \right), \quad (6)$$

where  $S_{\min}$  and  $S_{\max}$  represent the sets of the  $k$  least- and most-frequently selected prompts, respectively.  $\lambda_{\text{freq}}$  is a weighting factor. This regularization encourages underutilized prompts to participate more in learning and discourages prompt collapse into a narrow set of overused patterns.

2) **Prompt Knowledge Orthogonalization:** To differentiate the knowledge encoded by prompt pairs, we apply an orthogonalization constraint over the query and value prompts in the pool. By minimizing the cosine similarity between prompt embeddings, we reduce redundancy and promote representation diversity. The orthogonalization

**Algorithm 1** EgoPrompt Training Procedure

---

**Require:** Component-specific prompts  $\mathbf{p}_v, \mathbf{p}_n$ ; number of iterations  $T_1, T_2$  for two training stages; Unified Prompt Pool  $\mathcal{P}$  consisting of  $P$  query-value prompt pairs  $\{(\mathbf{q}^P, \mathbf{v}^P)\}_{P=1}^P$ ; Projector  $\text{Proj}(\cdot)$ .

- 1: // **Stage 1: Component-Specific Prompt Learning**
- 2: **for**  $t = 1$  to  $T_1$  **do**
- 3:   Extract component-specific video features  $f_v, f_n$ .
- 4:   Obtain the classification loss  $\mathcal{L}_{ce}$  and knowledge-guided  $\mathcal{L}_{kg}$  loss for both components via Eq. 1 and Eq. 2.
- 5:   Optimize  $\mathbf{p}_v$  and  $\mathbf{p}_n$ .
- 6: **end for**
- 7: // **Stage 2: Implicit Pattern Interaction Learning**
- 8: **for**  $t = 1$  to  $T_2$  **do**
- 9:   Extract component-specific features  $f_v, f_n$ .
- 10:   Retrieve top- $k$  query prompts  $\mathbf{Q}_v, \mathbf{Q}_n$  and their corresponding value prompts  $\mathbf{V}_v, \mathbf{V}_n$  from pool  $\mathcal{P}$  via Eq. 3.
- 11:   Compute fused representation  $f_s$  via Eq. 4 and Eq. 5.
- 12:   Obtain the unified space construction objective  $\mathcal{L}_{uni}$  via Eq. 8.
- 13:   Optimize the Unified Prompt Pool  $\mathcal{P}$  and  $\text{Proj}(\cdot)$ .
- 14: **end for**
- 15: **return**  $\mathbf{p}_v, \mathbf{p}_n, \mathcal{P}$ , and  $\text{Proj}(\cdot)$ .

---

loss is defined as:

$$\mathcal{L}_{orth} = \frac{1}{P(P-1)} \sum_{i=1}^P \sum_{j=1, j \neq i}^P \left( \left| \frac{\mathbf{q}_i \cdot \mathbf{q}_j}{\|\mathbf{q}_i\| \|\mathbf{q}_j\|} \right| + \left| \frac{\mathbf{v}_i \cdot \mathbf{v}_j}{\|\mathbf{v}_i\| \|\mathbf{v}_j\|} \right| \right), \quad (7)$$

where  $|\cdot|$  denotes the absolute value operator,  $\mathbf{q}_i \cdot \mathbf{q}_j$  and  $\mathbf{v}_i \cdot \mathbf{v}_j$  are the dot product between the  $i$ -th and  $j$ -th query and value prompts. This regularization  $\mathcal{L}_{orth}$  encourages prompt pairs to occupy distinct subspaces, thereby enhancing the diversity of captured semantic patterns. Overall, the final training objective for unified pattern space construction combines the cross-entropy losses from both verb and noun branches with the two regularization terms:

$$\mathcal{L}_{uni} = \mathcal{L}_{ce}^v + \mathcal{L}_{ce}^n + \lambda_{freq} \mathcal{L}_{freq} + \lambda_{orth} \mathcal{L}_{orth}, \quad (8)$$

where  $\lambda_{freq}$  and  $\lambda_{orth}$  are balancing coefficients controlling the strength of the frequency and orthogonality regularizations, respectively. We provide hyperparameter analysis in the supplementary materials. By enforcing both usage balance and semantic orthogonality, the **Diverse Pool Criteria** ensures that the unified prompt space remains rich, informative, and well-structured, supporting robust egocentric action recognition. We summarize the complete training procedure in Algorithm 1.

## 4 Experiments

### 4.1 Datasets

**Ego4D** [10]. We follow X-MIC [29] and use a subset of the Ego4D Forecasting Hands and Object (FHO) benchmark. It is annotated with fine-grained noun and verb labels. The training set consists of 64K video clips, with 521 noun classes and 117 verb classes. The test set includes 33K clips with an average clip duration of 8 seconds, amounting to approximately 215 hours of video, excluding irrelevant background clips.

**Epic-Kitchens** [11] includes 67K video clips for training and 10K video clips for validation and testing. Each clip is 3.5 seconds on average, amounting to approximately 70 hours of video (irrelevant

background clips are excluded). The dataset focuses on kitchen activities and is annotated for 300 noun classes and 97 verb classes. **EGTEA** [9] is only used for testing generalization performance. Its training set includes 8,000 video clips with an average duration of 3.2 seconds, while the test split comprises 6,000 video clips. The dataset provides 20 verb classes and 54 noun classes.

### 4.2 Experimental Setup

**Training Details.** Our implementation is adapted from the public code of the X-MIC [29]. To make a fair comparison, we adopt dual-encoder architecture based on the LAVILA [8] with **Timesformer-Large** as our default backbone for all algorithms. The prompt’s length is set to 4 for both video, textual modalities, and prompt pool construction. For initialization, we employ hand-crafted templates in the form of “a video of a [CLASS] action” for verb classes and “a video of actioning on [CLASS]” for noun classes. The Adamw [38] optimizer is applied for optimization with the learning rate of  $1e-4$  for formal training and  $2e-5$  for warmup, and the batch size in Ego4D and Epic-Kitchen is set for 64 and 32, and training epochs is 5, including 3 warm-up epochs. All experiments are conducted on 4×A800 GPUs with a single training run lasting 8 hours.

**Baselines.** To introduce the prompt-learning technique to egocentric action recognition, we re-implement several recent CoOp-based works, e.g., MaPLe [16] and KgCoOp [25]. In addition, we introduce some state-of-the-art egocentric action recognition methods, X-MIC [29], AoP [6], and POV [7]. To make a comprehensive comparison, we also include results from conventional supervised fine-tuning baselines, as illustrated in Fig. 3 and Fig. 4. All these methods are evaluated under consistent protocols on the following proposed benchmarks.

**Benchmarks.** We evaluate EgoPrompt under two standard egocentric action recognition benchmarks: 1) *Within- and Cross-dataset generalization*: In this setting, models are trained on a single dataset and evaluated on both the training domain (within-dataset) and an unseen testing domain (cross-dataset). This benchmark assesses the model’s ability to not only fit the source domain but also to generalize to new environments. Accordingly, we report performance on both domains to reflect the model’s adaptability to distribution shifts. 2) *Base-to-novel class generalization*: Here, datasets are divided into base and novel subsets based on whether the data label is shared across training and testing datasets or appears exclusively in the testing dataset. The model is trained solely on the base class subset of the training dataset and evaluated on both the base and novel classes within the testing dataset. Performance on base classes indicates the model’s ability to transfer category-level knowledge across domains, while performance on novel classes reflects its zero-shot recognition capability.

**Evaluation Metrics.** To assess generalization fairly, we adopt two complementary metrics: *Average Accuracy*: Measures the overall accuracy across all test samples, capturing instance-level correctness. *Class Average Accuracy*: Computes the mean accuracy across classes, mitigating bias from class imbalance and emphasizing performance consistency on both frequent and rare classes. Together, these metrics provide a holistic view of the model’s generalization behavior from both sample-level and class-level perspectives.



**Table 1: Comparison on the within- and cross-dataset generalization setting. The superscript of  $\oplus$  denotes the CoOp-based and  $\odot$  denotes the egocentric-based algorithms., “hm” is short for harmonic average.**

	Trained on Ego4D (E4D)						Trained on Epic-Kitchens (EK)					
	Nouns			Verbs			Nouns			Verbs		
	E4D	EK	hm	E4D	EK	hm	E4D	EK	hm	E4D	EK	hm
CoOp <sup><math>\oplus</math></sup> [19]	31.46	33.20	32.31	22.50	23.08	22.79	14.80	34.70	20.75	15.42	53.17	23.91
CoCoOp <sup><math>\oplus</math></sup> [20]	35.71	33.06	34.33	24.26	29.86	26.77	12.76	36.93	18.97	14.28	55.20	22.69
CLIP-Adapter <sup><math>\oplus</math></sup> [18]	33.20	29.80	31.41	23.68	31.92	27.19	11.40	37.85	17.52	16.13	56.77	25.12
A5 <sup><math>\oplus</math></sup> [35]	34.58	31.42	32.92	24.80	27.65	26.15	13.96	39.73	20.66	16.79	56.05	25.84
Vita-CLIP <sup><math>\oplus</math></sup> [36]	33.68	30.10	31.79	23.46	28.90	25.90	16.47	40.08	23.35	17.30	55.42	26.37
POV <sup><math>\odot</math></sup> [7]	37.60	32.86	35.07	27.46	34.80	30.70	13.75	37.70	20.15	16.74	53.08	25.45
X-MIC <sup><math>\odot</math></sup> [29]	35.85	28.26	31.61	28.27	39.49	32.95	11.45	44.07	18.17	16.01	53.02	24.60
KgCoOp <sup><math>\oplus</math></sup> [25]	36.71	34.18	35.40	22.40	43.17	29.50	17.28	41.72	24.44	18.93	52.70	27.85
AoP <sup><math>\odot</math></sup> [6]	34.60	34.26	34.43	24.63	39.76	30.42	11.62	38.76	17.88	13.16	49.45	20.79
MaPLe <sup><math>\oplus</math></sup> [16]	39.87	32.82	36.00	25.53	44.24	32.38	17.31	41.62	24.45	18.26	58.10	27.79
EgoPrompt	<b>42.93</b>	<b>35.75</b>	<b>39.01</b>	<b>29.71</b>	<b>47.89</b>	<b>36.67</b>	<b>19.45</b>	<b>44.58</b>	<b>27.08</b>	<b>20.78</b>	<b>61.40</b>	<b>31.05</b>

**Table 2: Comparison on the base-to-novel class generalization setting. The experiment results are pre-trained on Ego4D and evaluated on Epic-Kitchen.**

Pre-trained dataset: E4D Evaluation dataset: EK	Average Accuracy						Class Average Accuracy					
	Nouns			Verbs			Nouns			Verbs		
	base	novel	hm	base	novel	hm	base	novel	hm	base	novel	hm
CoOp <sup><math>\oplus</math></sup>	34.16	23.18	27.62	24.27	2.04	3.76	20.05	5.18	8.23	10.07	1.40	2.46
CoCoOp <sup><math>\oplus</math></sup>	35.17	23.07	27.86	24.63	1.52	2.86	20.98	5.71	8.98	11.63	1.71	2.98
CLIP-Adapter <sup><math>\oplus</math></sup>	33.28	21.72	26.29	30.74	2.71	4.98	19.76	4.98	7.96	11.92	2.04	3.48
A5 <sup><math>\oplus</math></sup>	34.10	22.90	27.40	27.60	2.80	5.08	20.73	3.48	5.96	10.37	1.82	3.10
Vita-CLIP <sup><math>\oplus</math></sup>	32.18	22.67	26.60	29.14	2.63	4.82	19.64	3.70	6.23	11.31	1.62	2.83
POV <sup><math>\odot</math></sup>	34.95	22.17	27.13	33.85	2.09	3.94	20.85	4.07	6.81	11.68	1.79	3.10
X-MIC <sup><math>\odot</math></sup>	34.32	23.00	27.54	30.17	2.42	4.48	20.46	5.02	8.06	10.80	1.93	3.27
AoP <sup><math>\odot</math></sup>	36.33	24.01	28.91	41.52	2.13	4.05	21.61	6.80	10.34	12.30	2.60	4.29
MaPLe <sup><math>\oplus</math></sup>	35.70	24.13	28.80	46.09	2.68	5.07	21.13	6.32	9.73	12.76	2.15	3.68
EgoPrompt	<b>36.91</b>	<b>24.34</b>	<b>29.34</b>	<b>50.70</b>	<b>3.07</b>	<b>5.79</b>	<b>21.83</b>	<b>7.14</b>	<b>10.76</b>	<b>13.77</b>	<b>3.21</b>	<b>5.21</b>

### 4.3 Within- and Cross-Dataset Generalization

**Prompt Learning vs. Fine-tuning vs. Zero-shot.** As shown in Fig. 3, prompt-based approaches (e.g., X-MIC and EgoPrompt) significantly outperform both zero-shot CLIP and fully fine-tuned (SFT) baselines in cross-dataset generalization. It is because zero-shot models suffer from poor alignment to downstream distributions, and SFT tends to overfit to the source domain, while prompt learning can help achieve a better balance between adaptability and generalization. Notably, EgoPrompt achieves the highest performance across both datasets and categories, especially in verb classification. For example, on Epic-Kitchens, EgoPrompt improves verb accuracy from 5.3% (zero-shot) and 16.8% (SFT) to 47.9%, clearly demonstrating the effectiveness of structured prompt learning in capturing transferable action semantics. Compared to X-MIC, which also introduces learnable modules, EgoPrompt further enhances generalization by modeling cross-component interactions and latent semantics, leading to consistent gains on both nouns and verbs.

**Ego4D & EpicKitchen.** We present a detailed comparison of generalization performance in Table 1, where models are trained on the source dataset and tested on both the source (i.e., within-dataset) and target datasets (i.e., cross-dataset) for noun and verb classification. (1) Within-dataset performance: EgoPrompt achieves the highest within-dataset accuracy on E4D (42.93% for nouns and 29.71% for verbs) and EK (44.58% for nouns and 61.40% for verbs). These results exceed the prior best-performing method MaPLe by +3.06% and +4.18% in E4D and +2.96% and +3.30% in EK, showing that our method not only generalizes well but also retains strong fitting ability on the within dataset. (2) Cross-dataset performance: EgoPrompt also outperforms all baselines on the unseen test set, reaching 35.75% noun and 47.89% verb in EK testing and 19.45%

noun and 20.78% verb in E4D testing. These consistent improvements are also reflected in the hm, with EgoPrompt outperforming MaPLe by +3.65% nouns and +4.29% verbs in EK testing and +2.14% nouns and 2.52% verbs in E4D testing. As a result, the experimental results demonstrate the effectiveness of EgoPrompt in both within- and cross-domain generalization gains.

**Ego4D & EGTEA.** To further evaluate the cross-dataset generalization capability of EgoPrompt, we report its performance on the EGTEA dataset in Table 3, where models are trained on Ego4D and tested directly on EGTEA without any adaptation. EgoPrompt achieves the best performance across all metrics, reaching 32.8% on noun classification and 40.3% on verb classification, surpassing all baselines by a clear margin. Compared to MaPLe, which previously achieved 30.7% (nouns) and 36.2% (verbs), EgoPrompt improves performance by +2.1% and +4.1%, respectively. These improvements are further reflected in the hm, where EgoPrompt achieves 36.2%, outperforming MaPLe (33.2%) and other methods, including X-MIC (30.3%) and KgCoOp (32.0%). Together with the results on Epic-Kitchens, these findings demonstrate that EgoPrompt not only retains strong within-domain and base class performance but also generalizes robustly to unseen domains and novel classes, effectively bridging the gap across egocentric video datasets.

### 4.4 Base-to-Novel Class Generalization

In Table 2, we evaluate the generalization performance of various methods under the base-to-novel setting, where models are trained on base classes in Ego4D and tested on both base and novel classes in Epic-Kitchens. EgoPrompt demonstrates the strongest performance across all metrics. It achieves substantial gains on base

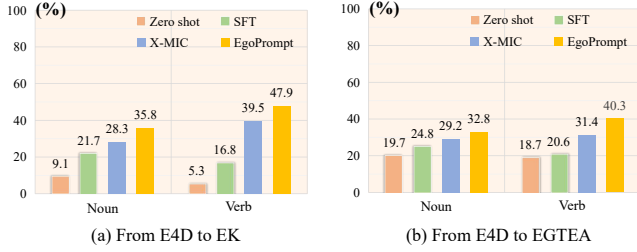


Figure 3: Comparison of the generalization performance. The sub-caption denotes the training and testing dataset in this cross-dataset generalization setting

Table 3: Cross-dataset generalization on EGTEA.

	E4D (within)			EGTEA (cross)		
	nouns	verbs	hm	nouns	verbs	hm
X-MIC	35.9	29.0	32.5	29.2	31.4	30.3
KgCoOp	36.7	22.4	27.8	29.8	34.6	32.0
AoP	34.6	24.6	28.8	30.1	32.6	31.3
MaPLE	39.9	25.5	31.1	30.7	36.2	33.2
EgoPrompt	<b>42.9</b>	<b>29.7</b>	<b>35.1</b>	<b>32.8</b>	<b>40.3</b>	<b>36.2</b>

class recognition, particularly for verbs, reaching 50.70% in average accuracy—surpassing CoOp (24.27%) and MaPLE (46.09%) by a large margin. This highlights EgoPrompt’s ability to extract more transferable action semantics in commonly occurring categories. In novel classes, all methods exhibit limited performance due to the inherent challenge of few-shot generalization. EgoPrompt slightly improves novel verb accuracy to 3.07%, compared to 2.04% of CoOp and 2.68% of MaPLE. Although this marks the highest result among all baselines, the absolute value remains low, indicating that the egocentric model still struggles in long-tail scenarios with scarce training data. From a class-level perspective (Class Average Accuracy), EgoPrompt again leads on both nouns and verbs across base and novel splits. Notably, the hm for verbs improves to 5.21%, which is +1.53% higher than MaPLE, confirming the consistent gains of our method. In summary, EgoPrompt excels in modeling base class semantics and offers modest gains on novel classes.

#### 4.5 Ablation Study

EgoPrompt combines prompt learning with the characteristics of egocentric action recognition tasks. Therefore, we conducted a series of ablation studies to validate the effectiveness of the EgoPrompt. All of our ablation experiments are pre-trained on Ego4D. **Effect of EgoPrompt training stage design.** EgoPrompt adopts a two-stage training strategy to balance component-specific representation learning and cross-component interaction. Specifically, Stage 1 focuses on Dual-Branch Prompt Learning, while Stage 2 captures implicit pattern interactions through the Unified Prompt Pool. To evaluate the contribution of each stage and the training strategy, we conduct an ablation study as shown in Table 5. We consider four variants: training with only Stage 1, only Stage 2, joint training of Stage 1 and 2 from scratch (denoted as “Stage 1+2”), and the full two-stage training used in EgoPrompt. Results show that removing either stage leads to notable performance degradation. Training with only Stage 1 yields 31.28% (nouns) and 41.92% (verbs)

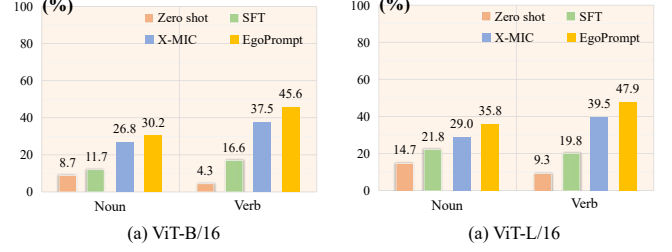


Figure 4: Adaptability of EgoPrompt on different backbones. The above results are collected from cross-dataset generalization, “From E4D to EK” setting.

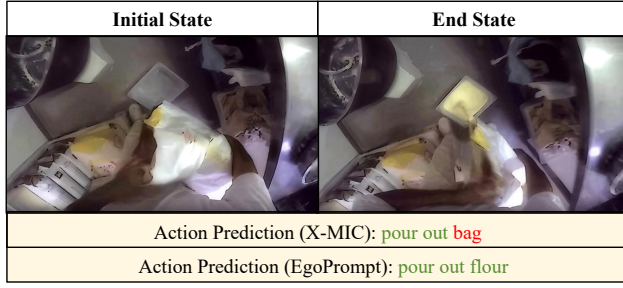
Table 4: Effect of modules in EgoPrompt.

Stage 1	Stage 2	Stage 1+2	E4D (within)		EK (cross)	
			nouns	verbs	nouns	verbs
✓	-	-	42.17	<b>29.82</b>	31.28	41.92
-	✓	-	40.18	28.31	33.15	46.30
-	-	✓	40.65	27.90	32.18	44.20
✓	✓	-	<b>42.93</b>	<b>29.71</b>	<b>35.75</b>	<b>47.89</b>

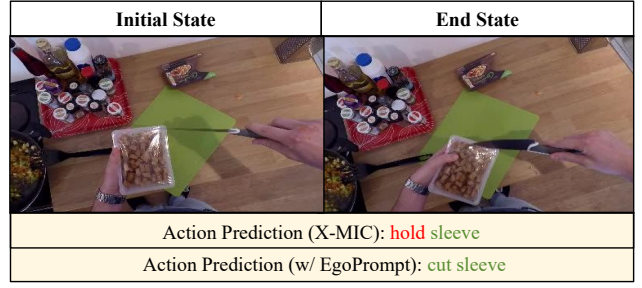
on Epic-Kitchens, while only Stage 2 gives slightly higher scores at 33.15% and 46.30%, indicating that interaction modeling contributes more to cross-domain generalization. Interestingly, jointly training both stages from scratch (“Stage 1+2”) fails to reach the same performance as our two-stage strategy, dropping to 32.18% (nouns) and 44.20% (verbs). This confirms that sequential optimization—first learning disentangled component-specific prompts and then modeling their interaction—is more effective for robust and transferable egocentric representations. These results validate the necessity of the proposed two-stage paradigm in EgoPrompt for achieving optimal performance.

**Adaptability across backbone scales.** To assess the scalability and adaptability of EgoPrompt under varying model capacities, we evaluate its performance on two video backbones: ViT-B/16 and ViT-L/16. The results are illustrated in Fig. 4. Across both backbones, EgoPrompt consistently outperforms zero-shot CLIP, supervised fine-tuning (SFT), and the adapter-based method X-MIC for both noun and verb recognition. On ViT-B/16, EgoPrompt achieves 30.2% (nouns) and 45.6% (verbs), exceeding X-MIC by +3.4% and +8.1%, respectively. When scaled up to ViT-L/16, the performance further improves to 35.8% (nouns) and 47.9% (verbs), showcasing enhanced generalization with increased backbone capacity. These results demonstrate that the proposed EgoPrompt is not only effective but also robust and adaptable across different model scales.

**Effect of Constraints in Diverse Pool Criteria.** Table 5 presents an ablation study on the two regularization terms of the Diverse Pool Criteria: Prompt Selection Frequency Regularization ( $\mathcal{L}_{freq}$ ) and Prompt Orthogonalization ( $\mathcal{L}_{orth}$ ). Two terms are designed to encourage a more balanced and diverse utilization of prompt pairs in the Unified Prompt Pool, respectively. Without either constraint, the model exhibits suboptimal generalization, particularly on cross-dataset (EK) evaluation, achieving only 34.46% in hm. Adding  $\mathcal{L}_{freq}$  alone leads to a clear improvement (hm +3.53% on EK), indicating that preventing over-reliance on a subset of prompt pairs helps improve verb classification in particular. Similarly, incorporating  $\mathcal{L}_{orth}$  alone also boosts performance (hm +5.43% on EK), especially



(a) Example 1: Noun Correction



(b) Example 2: Verb Correction

**Figure 5: Qualitative examples of EgoPrompt’s improvements.** (a) Noun correction: The baseline (X-MIC) incorrectly grounds the object to “bag,” while EgoPrompt leverages verb-centric semantics (e.g., deformable nature of “pour out”) to infer the correct noun “flour.” (b) Verb correction: The baseline misinterprets the action as “hold sleeve.” EgoPrompt captures the state change of the noun (i.e., sleeve being cut), which contradicts the static nature of “hold” and corrects the verb to “cut.” Green and red highlights indicate correct and incorrect predictions, respectively.

**Table 5: Effect of Components in Diverse Pool Criteria.**

$\mathcal{L}_{freq}$	$\mathcal{L}_{orth}$	E4D (within)			EK (cross)		
		nouns	verbs	hm	nouns	verbs	hm
-	-	38.31	24.71	30.04	30.17	40.16	34.46
✓	-	41.82	27.64	33.28	32.10	46.52	37.99
-	✓	40.96	28.54	33.64	34.60	47.08	39.89
✓	✓	<b>42.93</b>	29.71	35.12	<b>35.75</b>	<b>47.89</b>	40.94

on noun recognition, by encouraging prompt diversity and reducing redundancy within the pool. When both regularizations are applied jointly, the model achieves the highest performance across all metrics, with 35.12% hm on E4D and 40.94% on EK. These results confirm the complementary benefits of frequency balancing and orthogonal diversity, which together enhance the robustness and generalizability of the learned prompt space.

**Effect of Unified Prompt Pool size.** Table 6 investigates the impact of the prompt pool size ( $P$ ) on model performance. A smaller pool size (e.g.,  $P = 4$ ) provides limited pattern diversity, resulting in lower accuracy, particularly for nouns (hm = 33.88%). As the pool size increases, performance improves across both within- and cross-dataset settings. The best results are obtained when  $P = 16$ , yielding the highest harmonic mean for both noun (35.12%) and verb (40.94%) classification. Interestingly, further enlarging the pool to  $P = 32$  leads to a slight performance drop, suggesting that overly redundant prompts may dilute the selection quality and weaken discriminative learning. These findings indicate that a moderate prompt pool size strikes a favorable trade-off between pattern diversity and prompt effectiveness. In particular, setting  $P = 16$  enables the model to capture rich yet manageable semantic patterns for more robust cross-domain generalization. More Results can be found in the supplementary materials.

#### 4.6 Qualitative Analysis

We present two representative examples in Figure 5 to highlight how EgoPrompt enhances the model’s ability to reason over separate verb-noun semantics in egocentric videos. In **Example 1** (Fig. 5 (a)), the baseline model mispredicts the action as “pour out bag,” mistakenly linking the acted object to the “bag.” EgoPrompt

**Table 6: Effect of Prompt Pool Size.**

Pool Size $P$	Nouns			Verbs		
	E4D	EK	hm	E4D	EK	hm
4	41.60	28.58	33.88	<b>36.14</b>	46.50	40.67
8	42.32	<b>30.00</b>	35.11	34.98	47.20	40.18
<b>16</b>	<b>42.93</b>	<b>29.71</b>	<b>35.12</b>	<b>35.75</b>	<b>47.89</b>	<b>40.94</b>
32	42.70	29.11	34.62	33.26	45.84	38.55

captures the verb-centric semantic attribute of “pour out,” which typically involves a deformable or flowable object, successfully excluding the “bag” as a candidate. Based on this, EgoPrompt corrects the prediction to the right answer. In **Example 2** (Fig. 5 (b)), the baseline incorrectly labels the action as “hold sleeve,” overlooking the dynamic interaction occurring in the scene. In contrast, EgoPrompt leverages the object state change of the “sleeve” (being cut off), which is encoded in the noun representation. The cues contained in noun representation conflict with the semantics of “hold,” which implies a stable and static object state. By utilizing this inconsistency, EgoPrompt effectively rules out “hold” as a plausible verb and instead predicts “cut sleeve,” which better aligns with the observed scene dynamics.

## 5 Conclusion

In this paper, we propose EgoPrompt, a novel prompt learning framework tailored for egocentric action recognition. EgoPrompt explicitly models the component interaction through a Unified Prompt Pool, enabling the construction of a more semantically aligned and context-aware fused representation. To further enhance the effectiveness and diversity of learned prompt pairs, we introduce the Diverse Pool Criteria, which encourages balanced usage and orthogonal semantics across prompts. Extensive experiments demonstrate the superior performance of EgoPrompt.

**Discussion.** While EgoPrompt improves generalization and interpretability in egocentric recognition, challenges such as the long-tail distribution remain. Future research could explore incorporating temporal attention or dynamic prompt adaptation strategies, extending the temporal HOI understanding capability of EgoPrompt.



## 6 Acknowledgment

This work was supported by the National Natural Science Foundation of China under Grants 62036012, U23A20387, 62322212, in part by the Pengcheng Laboratory Research Project under Grant PCL2023A08, and also in part by the Postdoctoral Fellowship Program of CPSF under Grant Number GZC20251036.

## References

- [1] Feng Li, Renrui Zhang, Hao Zhang, Yuanhan Zhang, Bo Li, Wei Li, Zejun Ma, and Chunyuan Li. 2024. Llava-next-interleave: Tackling multi-image, video, and 3d in large multimodal models. *arXiv preprint arXiv:2407.07895* (2024).
- [2] Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Yanwei Li, Ziwei Liu, and Chunyuan Li. 2024. Llava-onevision: Easy visual task transfer. *arXiv preprint arXiv:2408.03326* (2024).
- [3] Shuhan Tan, Tushar Nagarajan, and Kristen Grauman. 2023. Egodistill: Egocentric head motion distillation for efficient video understanding. *Advances in Neural Information Processing Systems* 36 (2023), 33485–33498.
- [4] Tsukasa Shiota, Motohiro Takagi, Kaori Kumagai, Hitoshi Seshimo, and Yushi Aono. 2024. Egocentric action recognition by capturing hand-object contact and object state. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*. 6541–6551.
- [5] Chuhan Zhang, Ankush Gupta, and Andrew Zisserman. 2023. Helping hands: An object-aware ego-centric video recognition model. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 13901–13912.
- [6] Dibyadip Chatterjee, Fadime Sener, Shugao Ma, and Angela Yao. 2024. Opening the vocabulary of egocentric actions. *Advances in Neural Information Processing Systems* 36 (2024).
- [7] Boshen Xu, Sipeng Zheng, and Qin Jin. 2023. POV: Prompt-Oriented View-Agnostic Learning for Egocentric Hand-Object Interaction in the Multi-View World. In *Proceedings of the 31st ACM International Conference on Multimedia*. 2807–2816.
- [8] Yue Zhao, Ishan Misra, Philipp Krähenbühl, and Rohit Girdhar. 2023. Learning video representations from large language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 6586–6597.
- [9] Yin Li, Miao Liu, and James M Reh. 2018. In the eye of beholder: Joint learning of gaze and actions in first person video. In *Proceedings of the European conference on computer vision (ECCV)*. 619–635.
- [10] Kristen Grauman, Andrew Westbury, Eugene Byrne, Zachary Chavis, Antonino Furnari, Rohit Girdhar, Jackson Hamburger, Hao Jiang, Miao Liu, Xingyu Liu, et al. 2022. Ego4d: Around the world in 3,000 hours of egocentric video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 18995–19012.
- [11] Dima Damen, Hazel Doughty, Giovanni Maria Farinella, Antonino Furnari, Evangelos Kazakos, Jian Ma, Davide Moltisanti, Jonathan Munro, Toby Perrett, Will Price, et al. 2022. Rescaling egocentric vision: Collection, pipeline and challenges for epic-kitchens-100. *International Journal of Computer Vision* (2022), 1–23.
- [12] Yifei Huang, Guo Chen, Jilan Xu, Mingfang Zhang, Lijin Yang, Baoqi Pei, Hongjie Zhang, Dong Lu, Yali Wang, Limin Wang, and Yu Qiao. 2024. EgoExoLearn: A Dataset for Bridging Asynchronous Ego- and Exo-centric View of Procedural Activities in Real World. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*.
- [13] Chaofan Chen, Xiaoshan Yang, and Changsheng Xu. 2025. Pseudo Informative Episode Construction for Few-Shot Class-Incremental Learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 39. 15749–15757.
- [14] Menglin Jia, Luming Tang, Bor-Chun Chen, Claire Cardie, Serge Belongie, Bharath Hariharan, and Ser-Nam Lim. 2022. Visual prompt tuning. In *European Conference on Computer Vision*. Springer, 709–727.
- [15] Hantao Yao, Rui Zhang, and Changsheng Xu. 2024. TCP: Textual-based Class-aware Prompt tuning for Visual-Language Model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 23438–23448.
- [16] Muhammad Uzair Khattak, Hanoona Rasheed, Muhammad Maaz, Salman Khan, and Fahad Shahbaz Khan. 2023. MaPL: Multi-modal Prompt Learning. *arXiv:2210.03117* [cs.CV].
- [17] Zangwei Zheng, Xiangyu Yue, Kai Wang, and Yang You. 2022. Prompt Vision Transformer for Domain Generalization. *arXiv:2208.08914* [cs.CV].
- [18] Peng Gao, Shijie Geng, Renrui Zhang, Teli Ma, Rongyao Fang, Yongfeng Zhang, Hongsheng Li, and Yu Qiao. 2024. Clip-adapter: Better vision-language models with feature adapters. *International Journal of Computer Vision* 132, 2 (2024), 581–595.
- [19] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. 2022. Learning to Prompt for Vision-Language Models. *International Journal of Computer Vision* 130, 9 (July 2022), 2337–2348. doi:10.1007/s11263-022-01653-1
- [20] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. 2022. Conditional Prompt Learning for Vision-Language Models. *arXiv:2203.05557* [cs.CV].
- [21] Hantao Yao, Rui Zhang, Huaihai Lyu, Yongdong Zhang, and Changsheng Xu. 2025. Bi-Modality Individual-Aware Prompt Tuning for Visual-Language Model. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 47, 8 (2025), 6352–6368.
- [22] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In *International conference on machine learning*. PMLR, 8748–8763.
- [23] Chaofan Chen, Xiaoshan Yang, Jinpeng Zhang, Bo Dong, and Changsheng Xu. 2023. Category knowledge-guided parameter calibration for few-shot object detection. *IEEE Transactions on Image Processing* 32 (2023), 1092–1107.
- [24] Muhammad Uzair Khattak, Syed Talal Wasim, Muzammal Naseer, Salman Khan, Ming-Hsuan Yang, and Fahad Shahbaz Khan. 2023. Self-regulating prompts: Foundational model adaptation without forgetting. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 15190–15200.
- [25] Hantao Yao, Rui Zhang, and Changsheng Xu. 2023. Visual-language prompt tuning with knowledge-guided context optimization. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 6757–6767.
- [26] Zifeng Wang, Zizhao Zhang, Chen-Yu Lee, Han Zhang, Ruoxi Sun, Xiaoqi Ren, Guolong Su, Vincent Perot, Jennifer Dy, and Tomas Pfister. 2022. Learning to prompt for continual learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 139–149.
- [27] Shraman Pramanick, Yale Song, Sayan Nag, Kevin Qinghong Lin, Hardik Shah, Mike Zheng Shou, Rama Chellappa, and Pengchuan Zhang. 2023. EgoVLPv2: Egocentric Video-Language Pre-training with Fusion in the Backbone. *arXiv:2307.05463* [cs.CV] <https://arxiv.org/abs/2307.05463>
- [28] Sijie Cheng, Zhicheng Guo, Jingwen Wu, Kechen Fang, Peng Li, Huaping Liu, and Yang Liu. 2024. EgoThink: Evaluating First-Person Perspective Thinking Capability of Vision-Language Models. *arXiv:2311.15596* [cs.CV] <https://arxiv.org/abs/2311.15596>
- [29] Anna Kukleva, Fadime Sener, Edoardo Remelli, Bugra Tekin, Eric Sauser, Bernt Schiele, and Shugao Ma. 2024. X-MIC: Cross-Modal Instance Conditioning for Egocentric Action Generalization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 26364–26373.
- [30] Kevin Qinghong Lin, Alex Jinpeng Wang, Mattia Soldan, Michael Wray, Rui Yan, Eric Zhongcong Xu, Difei Gao, Rongcheng Tu, Wenzhe Zhao, Weijie Kong, Chengfei Cai, Hongfa Wang, Dima Damen, Bernard Ghanem, Wei Liu, and Mike Zheng Shou. 2022. Egocentric Video-Language Pretraining. *arXiv:2206.01670* [cs.CV] <https://arxiv.org/abs/2206.01670>
- [31] Yuheng Ji, Huajie Tan, Jiayu Shi, Xiaoshuai Hao, Yuan Zhang, Hengyuan Zhang, Pengwei Wang, Mengdi Zhao, Yao Mu, Pengju An, Xinda Xue, Qinghang Su, Huaihai Lyu, Xiaolong Zheng, Jiaming Liu, Zhongyuan Wang, and Shanghang Zhang. 2025. RoboBrain: A Unified Brain Model for Robotic Manipulation from Abstract to Concrete. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 1724–1734.
- [32] BAAI RoboBrain Team, Mingyu Cao, Huajie Tan, Yuheng Ji, Minglan Lin, Zhiyu Li, Zhou Cao, Pengwei Wang, Enshen Zhou, Yi Han, Yingbo Tang, Xiangqi Xu, Wei Guo, Yaoxu Lyu, Yijie Xu, Jiayu Shi, Mengfei Du, Cheng Chi, Mengdi Zhao, Xiaoshuai Hao, Junkai Zhao, Xiaojie Zhang, Shanyu Rong, Huaihai Lyu, Zhengliang Cai, Yankai Fu, Ning Chen, Bolun Zhang, Lingfeng Zhang, Shuyi Zhang, Dong Liu, Xi Feng, Songjing Wang, Xiaodan Liu, Yance Jiao, Mengsi Lyu, Zhuo Chen, Chenrui He, Yulong Ao, Xue Sun, Zheqi He, Jingshu Zheng, Xi Yang, Donghai Shi, Kunchang Xie, Bochao Zhang, Shaokai Nie, Chunlei Men, Yonghua Lin, Zhongyuan Wang, Tiejun Huang, and Shanghang Zhang. 2025. RoboBrain 2.0 Technical Report. *arXiv:2507.02029* [cs.RO] <https://arxiv.org/abs/2507.02029>
- [33] Huiyu Wang, Mitesh Kumar Singh, and Lorenzo Torresani. 2023. Ego-only: Egocentric action detection without exocentric transferring. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 5250–5261.
- [34] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. 2022. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 16000–16009.
- [35] Chen Ju, Tengda Han, Kunhao Zheng, Ya Zhang, and Weidi Xie. 2022. Prompting visual-language models for efficient video understanding. In *European Conference on Computer Vision*. Springer, 105–124.
- [36] Syed Talal Wasim, Muzammal Naseer, Salman Khan, Fahad Shahbaz Khan, and Mubarak Shah. 2023. Vita-clip: Video and text adaptive clip via multimodal prompting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 23034–23044.
- [37] Chaofan Chen, Xiaoshan Yang, Changsheng Xu, Xuhui Huang, and Zhe Ma. 2021. Eckpn: Explicit class knowledge propagation network for transductive few-shot learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 6596–6605.
- [38] Ilya Loshchilov and Frank Hutter. 2017. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101* (2017).