SPATIO-TEMPORAL DECOUPLED KNOWLEDGE COM PENSATOR FOR FEW-SHOT ACTION RECOGNITION

Anonymous authors

Paper under double-blind review

ABSTRACT

Few-Shot Action Recognition (FSAR) is a challenging task that requires recognizing novel action categories with a few labeled videos. Recent works typically apply semantically coarse category names as auxiliary contexts to guide the learning of discriminative visual features. However, such context provided by the action names is too limited to provide sufficient background knowledge for capturing novel spatial and temporal concepts in actions. In this paper, we propose **DIST**, an innovative **D**ecomposition-incorporation framework for FSAR that makes use of decoupled Spatial and Temporal knowledge provided by large language models to learn expressive multi-granularity prototypes. In the decomposition stage, we decouple vanilla action names into diverse spatio-temporal attribute descriptions (*i.e.*, action-related knowledge). Such commonsense knowledge complements semantic contexts from spatial and temporal perspectives. In the incorporation stage, we propose Spatial/Temporal Knowledge Compensators (SKC/PKC) to discover discriminative object- and frame-level prototypes, respectively. In SKC, objectlevel prototypes adaptively aggregate important patch tokens under the guidance of spatial knowledge. Moreover, in TKC, frame-level prototypes utilize temporal attributes to assist in inter-frame temporal relation modeling, further understanding diverse temporal patterns in videos. These learned prototypes at varying levels of granularity thus provide transparency in capturing fine-grained spatial details and dynamic temporal information, so as to enable accurate recognition of both appearance-centric and motion-centric actions. Experimental results show DIST achieves state-of-the-art results on four standard FSAR datasets (*i.e.*, Kinetics, UCF101, HMDB51 and SSv2-small). Full code will be released.

032 033

004

010 011

012

013

014

015

016

017

018

019

021

024

025

026

027

028

029

031

034 035

037

1 INTRODUCTION

With deep learning advancements [1, 2, 3, 4], significant progress has been made in the field of action recognition [5, 6, 7, 8, 9] recently. However, this success relies heavily on a large amount of manually-labeled samples, which are time-consuming and expensive to acquire. To alleviate the data-hunger issue, considerable works [10, 11, 12, 13, 14] have turned their attention to few-shot action recognition (FSAR), where there exist base action classes (seen) with a large volume of training examples, and novel action classes (unseen) with unlabeled samples. FSAR learns feature representation from base action classes, and then evaluates its generalization ability on novel action classes.

044 Modern FSAR solutions [15, 16, 17, 18, 19, 20, 21] are largely built upon metric-based metalearning paradigm [22], where the model learns class (prototype) representation and performs 046 prototype-query matching with respect to predefined or learned distance metrics. Among them, 047 the top-leading methods [11, 15, 23, 13] directly extract class-related spatio-temporal feature repre-048 sentation from raw visual signals. Though impressive, these methods lack a basic grasp of explicit action knowledge, struggling to learn new concepts in action classes, particularly under data-limited conditions. Recent works [24, 21, 25] transfer knowledge from pre-trained vision-language mod-051 els [26, 27] (e.g., CLIP [27]) to enhance FSAR model capability. However, these methods typically apply semantically coarse or ambiguous category names as auxiliary context information to com-052 pensate for visual features. Such context provided by the action names is too limited to provide enough background knowledge for video action understanding [28, 29, 30].

054 Spatial Knowledge \bigcirc LLM prompts bject 1: Containe Text 1. Given action label {class 056 Encod e}, please generate K mos bject K: Mouth \cap related objects 2. Given action label {class name}, please describe M states of each action in simple **Object-level** prototypes Inte ction Spatial metri and short words Visnal Prediction Visual features ╬ 060 Temporal metri 🕼 LLM 061 Interaction Temporal Kno 062 \bigcirc tep 1: 'Hold contain \bigcirc Drink Text 063 Frame-level prototypes 2: 'Bring containe Temporal knowledge to mouth' 064 \bigcirc **Decomposition Stage** Incorporation Stage ep M: Put container 065

066 Figure 1: Our main idea. Our approach decomposes category names into diverse spatio-temporal knowledge, and makes use of the decoupled knowledge to learn object-/frame- level prototypes, respectively.

068 To address this issue, we study how to effectively collect and leverage action-related commonsense 069 knowledge provided by large language models (LLMs), compensating for visual information and 070 thus enhancing the few-shot learning capacity. The LLM serves as a knowledge base to provide 071 action-related background commonsense descriptions from complementary spatial and temporal 072 perspectives, which we refer to as decoupled "spatio-temporal prior knowledge". Compared to 073 vanilla categories, spatio-temporal prior knowledge i) makes up for missing context information 074 to achieve semantic completeness, and ii) transforms unseen categories into known commonsense 075 descriptions, easily interpreted by pre-trained language models.

076 In light of the above, we develop a novel decomposition-incorporation framework for FSAR: **DIST**, 077 which firstly decomposes vanilla category names into diverse spatio-temporal attribute descriptions, and then incorporates decoupled commonsense knowledge and visual features to guide the learn-079 ing of object-level and frame-level prototypes. As shown in Fig. 1, for the **decomposition stage**, we make use of LLMs, to generate action-related commonsense descriptions (*i.e.*, external contex-081 tual information and different steps of action) from coarse category names. Such comprehensive descriptions complement semantic contexts from spatial and temporal perspectives. In the **incor**-083 poration stage, we propose Spatial/Temporal Knowledge Compensators (SKC/TKC), which incorporate decoupled prior knowledge and visual features to form discriminative object-level (spatial) 084 and frame-level (temporal) prototypes, respectively. Specifically, SKC aggregates important patches 085 into compact object prototypes by patch-level cross-attention within each frame, and further guides object-level prototype learning with the assistance of spatial prior knowledge. These object-level 087 prototypes filter out video noise and focus on informative image patches, which refer to the most 880 class-related ones and correspond to key entities. Meanwhile, TKC captures the temporal rela-089 tionships between frame-level prototypes through inter-frame interaction, and then enables these 090 prototypes to aggregate essential semantic information from temporal prior knowledge. The learned 091 different-level prototypes can capture fine-grained spatial details and dynamic temporal information, 092 respectively, so as to yield more accurate action recognition results.

093 Overall, our contributions are summarized as follows: • We pioneer the early exploration in mak-094 ing use of action-related prior knowledge for FSAR, and to achieve this end, we construct background commonsense descriptions provided by LLMs in a spatiotemporal-decoupled manner. 096 We propose a novel decomposition-incorporation framework that decouples category names into diverse spatial and temporal prior knowledge, and then incorporates them and visual features to learn 098 object-level and frame-level prototypes in a dual way. ⁽³⁾ We design Spatial/Temporal Knowledge Compensators (SKC/TKC) that inject decoupled prior knowledge into different-level prototypes to 099 capture fine-grained spatial details and dynamic temporal information. 100

101 To the best of our knowledge, we make the pioneering effort to explore the application of diverse 102 spatio-temporal prior knowledge in FSAR, aiming to effectively provide semantic contexts for the 103 learning of different-level prototypes from multiple perspectives. Different from simply combining 104 rich LLM-generated descriptions (*i.e.*, only one sentence about actions) and frame-level features in 105 a single branch, our framework conducts customized feature interaction for different branches to incorporate i) patch-level features and spatial knowledge; and ii) frame-level features and temporal 106 knowledge, respectively (see Table 5c in \$4.4). Such a framework can discover fine-grained spa-107 tial patterns and dynamic temporal patterns, hence producing more accurate few-shot results. To comprehensively evaluate our method, we conduct experiments on four gold-standard datasets (*i.e.*, HMDB51 [31], UCF101 [32], kinetics100 [33], and SSv2-small [34]). We empirically prove that DIST surpasses all existing state-of-the-arts and yields solid performance gains (1.7%-6.8% accuracy) under the 5-way 1-shot setting. Furthermore, we perform thorough ablation studies to dissect each component, both quantitatively and qualitatively. Our full implementation will be released.

113 114

2 RELATED WORK

115 116

117 Few-shot Image Classification. The objective of few-shot image classification [35, 36] is to recog-118 nize new categories with a small number of annotated samples. Existing methods can be roughly categorized into three groups: i) Augmentation-based methods [2, 37, 38] exploit various augmenta-119 tion strategies to alleviate the data scarcity dilemma, mainly including spatial deformation [39] and 120 feature augmentation [40, 41]; ii) Optimization-based methods [42, 43, 44, 45] learn optimization 121 states, like model initialization [42, 43] or step sizes [44, 45], to update models with a few gradient 122 steps; and iii) Metric-based methods [46, 47, 22, 48, 49, 50] learn a class representation (prototype) 123 by averaging embeddings belonging to the same class, and predict query (*i.e.*, test sample) labels 124 with respect to predefined [46, 47, 22, 48] or learned [49, 50] distance metric. Our work is more 125 closely related to the metric-based methods [47, 22], whereas we focus on few-shot action recog-126 nition – a more challenging task that requires handling videos encompassing a wealth of temporal 127 information due to common and distinct patterns in nearby frames [51]. 128

Few-shot Action Recognition (FSAR). FSAR is a challenging task with the goal of recognizing 129 previously unseen action classes (*i.e.*, query class) with a few labeled videos. Existing FSAR meth-130 ods [52, 53, 54, 13, 15] mainly belong to the metric-based meta-learning paradigm [47], which 131 learns class (prototype) representation and performs prototype-query matching based on the learned 132 distance metrics. These methods are mainly devoted to feature representation learning [15, 13] and 133 matching strategy exploration [52, 55, 56, 11, 18, 16, 13]. As a primary step, feature representation 134 *learning* helps models to learn expressive spatio-temporal features for further matching process. 135 Recent appoarches [15, 13] model temporal features through temporal attention operations [15] or 136 more detailed temporal-patch and temporal-channel interaction [13], and further exploit low-level 137 spatial features by patch-level information interaction within each frame or across frames [57]. Some others make use of video features in a whole task (*i.e.*, episode) to extract relevant discriminative 138 patterns [13, 16, 58] by a graph neural network [13] or attention relation modeling [16]. For *match*-139 ing strategy exploration, early works [52, 55, 56] aggregate the frame features into a single video 140 representation for video-level feature matching. Though straightforward, these methods suffer from 141 suboptimal performance due to neglecting the temporal cues in videos. To address this limitation, 142 the following approaches [11, 18, 16, 13, 23] devise various temporal alignment metrics for frame-143 level feature matching, e.g., frame-level alignment [11, 16], segment-level alignment [18], and even 144 frame-to-segment alignment [23] that is also common in realistic video matching. 145

Recent works [24, 59, 21] transfer knowledge from pre-trained vision-language models (*e.g.*, CLIP [27])toenhance FSAR model capability. Though promising, they heavily rely on semantically coarse or ambiguous category names as semantic source to provide action-related context. Such context is insufficient to offer enough background knowledge for video understanding. In contrast, our DIST represents the first effort in FSAR to decouple class names into diverse spatio-temporal attribute descriptions (*i.e.*, action-relevant knowledge) to complement semantic contexts. More significantly, such acquired decoupled knowledge is further injected into visual features to learn objectand frame-level prototypes in a dual way, so as to enable more accurate action recognition.

153 **Few-shot Learning with Semantic Information.** Recent works on few-shot learning [60, 61, 62] 154 integrate semantic information (provided by class labels) and visual information (extracted from 155 visual observations) to represent a novel class. Based on the levels at which modality information 156 fusion occurs, these methods can be roughly categorized into three groups: i) Prototype-level meth-157 ods [63, 64] model class (prototype) representation as a combination of visual prototypes and se-158 mantic prototypes obtained through word embeddings of class labels by attention mechanism [64] or 159 adaptive fusion mechanism [63]; ii) Classifier-level methods [60] enable classifiers to predict novel categories by incorporating auxiliary semantic information acquired from a graph convolutional 160 network [65]; and iii) Extractor-level methods [66, 62] consider semantic information as prompts to 161 tune the feature extractor, allowing the feature extractor to better focus on class-specific features.



Figure 2: Overview of DIST. Video inputs are first processed by the visual encoder of CLIP to obtain initial patch-level and frame-level features. Then, we leverage LLM to decompose vanilla action names into action-related background knowledge. Furthermore, SKC/TKC incorporate decoupled prior knowledge and visual features to form discriminative object- and frame-level prototypes for spatial/temporal matching. Finally, we can combine spatial and temporal matching results to obtain the merged query prediction.

Though impressive, they [63, 66, 62] typically directly apply semantically coarse category names as auxiliary information at different levels to compensate for visual features, lacking high-quality background knowledge to discover novel visual concepts, therefore struggling with adapting to unseen categories. Our contribution is orthogonal to previous studies, as we advance FSAR regime in the aspect of collecting and leveraging high-quality spatio-temporal prior knowledge. The LLM serves as a knowledge base to provide action-related background commonsense descriptions (*i.e.*, contextual information and different steps of actions). Then our work makes smart use of prior knowledge to reduce redundant visual features and enhance the semantic distinction of different class prototypes.

3 Method

3.1 PROBLEM FORMULATION

194 The goal of FSAR is to classify unlabeled test videos given a few (e.g., one or five) samples per class. 195 Under the few-shot setting, the model learns a feature representation on training classes \mathcal{D}_{base} and is 196 evaluated on testing classes \mathcal{D}_{novel} to emphasize its generalization ability on novel categories, where 197 $\mathcal{D}_{base} \cap \mathcal{D}_{novel} = \emptyset$. In the training stage, following [67, 16, 13], we train a few-shot learning model in an episodic way. Here, each episode (i.e., a standard M-way K-shot episode task) is formed by sampling M categories from \mathcal{D}_{base} . The M-way K-shot task consists of K labeled videos per 199 class as the support set \mathcal{S} , and a fraction of the rest samples from M classes as the query set \mathcal{Q} . The 200 episodic training for FSAR is achieved by minimizing, for each episode, the loss of the prediction on 201 samples in the query set, given the support set. In the inference stage, we randomly sample episode 202 tasks from \mathcal{D}_{novel} for evaluation, and report average results over multiple episode tasks. 203

204 205

190

191 192

193

3.2 OVERALL FRAMEWORK

206 We introduce DIST, which collects and leverages action-related commonsense knowledge provided 207 by LLMs to guide the learning of object- and frame-level prototypes for FSAR (Fig. 2). DIST con-208 sists of visual/text encoders and two knowledge compensators. The model takes the RGB frame 209 sequence of length T and corresponding action names as input. We first utilize the visual encoder 210 of CLIP [27] to get frame-level (*i.e.*, class token in each frame) and patch-level feature representa-tion, *i.e.*, $F \in \mathbb{R}^{T \times C}$ and $X \in \mathbb{R}^{T \times P \times C}$, where P is the number of tokens in each frame. Besides, 211 212 we prompt LLM with corresponding action names to generate decoupled spatial and temporal commonsense descriptions (§3.3). These descriptions are fed into frozen text encoder of CLIP to obtain spatial and temporal attribute features, *i.e.*, $Q_s \in \mathbb{R}^{G \times C}$ and $Q_t \in \mathbb{R}^{L \times C}$. We further design two 213 214 complementary modules to make use of decoupled spatio-temporal attributes: 1) Spatial Knowledge 215 Compensator (SKC) (§3.4) injects spatial attributes into patch-level features to explicitly learn compact object-level prototypes for object-level prototype matching; 2) Temporal Knowledge Compensator (TKC) (§3.5) incorporates temporal attributes and frame-level features to inform frame-level prototypes for frame-level prototype matching. Finally, we can combine spatial and temporal feature matching scores to obtain the merged query prediction.

220 221

222

3.3 DECOUPLED SPATIO-TEMPORAL ATTRIBUTE GENERATION

223 Spatial Attribute Generation. Naive category names provided limited commonsense knowledge to focus on action-related spatial contexts. Thus we make use of prior knowledge in LLMs [68, 69] 224 to generate detailed and informative spatial attributes for each category, *i.e.*, action-related object 225 instances and environment. Specifically, taking the action category "drink" as an example, to obtain 226 spatial attribute descriptions, we prompt ChatGPT [69] by "Given action label {drink}, please gen-227 erate $\{G\}$ most related objects for each class.", where G is empirically set to 6. This prompt returns 228 a set \mathcal{A} with G spatial attribute descriptions, such as "container; mouth; hand; ...". Then we encode 229 these spatial attributes via frozen CLIP text encoder to get spatial attribute features $Q_s \in \mathbb{R}^{G \times C}$. 230

Temporal Attribute Generation. Prior researches [24, 59] apply semantically coarse category 231 names as auxiliary information to guide temporal feature learning. However, such context provided 232 by action names is too limited to provide enough temporal context for action recognition. Thus, we 233 propose to utilize the abundant prior knowledge in LLMs [68, 69] to expand the coarse action names. 234 Temporal attribute descriptions generated by LLM are a collection of multiple atomic actions, which 235 describe the temporal evolution of an action. Concretely, taking the action category "drink" as an 236 example, to obtain temporal attribute descriptions, we prompt ChatGPT [69] by "Given action label 237 $\{drink\}, please describe \{L\}$ states of each action in simple and short words.", where we empirically 238 set L to 3. This prompt always returns a set \mathcal{B} with L temporal attribute descriptions, such as "Hold 239 container; Bring container to mouth; Put container; ...", which decompose one action class into multiple atomic actions in a step-by-step manner. Then we adopt the off-the-shelf text encoder of 240 CLIP to encode these descriptions and obtain temporal attribute features $Q_t \in \mathbb{R}^{L \times C}$. 241

242 243

3.4 SPATIAL KNOWLEDGE COMPENSATOR

244 Previous methods exploit spatial features by patch-level information interaction within each frame or 245 across frames [15, 57]. However, this leads to two issues: 1) Too many irrelevant patch tokens bring 246 redundant information, interfering with further spatial feature matching; 2) They fail to focus on 247 important objects without the guidance of spatial prior knowledge. Therefore, we investigate how to 248 better incorporate spatial attributes (§3.3) and patch-level visual features into compact object-level 249 prototypes to highlight potential target objects. To this end, as showcased in Fig. 3, proposed Spatial 250 Knowledge Compensator (SKC) summarizes discriminative spatial patterns via aggregating patch-251 level features into compact object-level prototypes (*i.e.*, patch aggregation), and further delivers the 252 union of spatial attribute knowledge and such object-level prototypes to enhance learned spatial 253 patterns (*i.e.*, attribute injection).

Patch Aggregation. We first introduce a set of learnable object-level prototypes to aggregate image content and highlight potential target objects. The prototypes are randomly initialized embeddings and represented as $P_o \in \mathbb{R}^{N \times C}$, where N is the number of object prototypes. Firstly, a selfattention layer is adopted for the N object prototypes to interact with each other in each frame. Then, these prototypes aim to adaptively aggregate action-related or object-related key patches in a sparse manner by patch-level cross-attention within each frame. Specifically, for patch tokens $X^l \in \mathbb{R}^{P \times C}$ in *l*-th frame, the process can be defined as:

261

268

$$\hat{\boldsymbol{P}}_{\mathrm{o}} = \mathrm{Softmax}(\boldsymbol{P}_{\mathrm{o}}\boldsymbol{K}_{\mathrm{p}}^{\top})\boldsymbol{V}_{\mathrm{p}} + \boldsymbol{P}_{\mathrm{o}},$$
(1)

where K_p and V_p are the linear transformation features of patch tokens X^l . This allows the objectlevel prototypes to capture discriminative spatial patterns.

Attribute Injection. To further encourage object prototypes to focus on action-related spatial context information, we deliver the union of spatial attribute knowledge and such object-level prototypes to discover fine-grained spatial patterns via the attention mechanism as follows:

$$\mathbf{P}_{\rm s} = \text{Softmax}(\hat{\mathbf{P}}_{\rm o} \mathbf{K}_{\rm q}^{\top}) \mathbf{V}_{\rm q} + \hat{\mathbf{P}}_{\rm o}, \tag{2}$$

where K_q and V_q are the linear transformation features of spatial attribute features Q_s , $P_s \in \mathbb{R}^{N \times C}$ is learned diverse object prototypes. Note that the operation of object prototypes in each frame is

270 271 272 Self-Attention Cross-Attention Cross-Attention 273 274 earnable object Object prototypes 275 prototypes × C K 276 Patch token Spatial attributes Patch Aggregation Attribute Injection 277

Figure 3: Illustration of Spatial Knowledge Compensator (SKC). SKC aims to learn discriminative objectlevel prototypes in a sparse aggregation manner via patch aggregation and attribute injection.

the same. By exchanging information for visual features and attribute features respectively, learned object-level prototypes filter out redundant information in videos and capture spatial details.

3.5 TEMPORAL KNOWLEDGE COMPENSATOR

How to incorporate temporal attribute features and frame-level features is essential for better FSAR performance since temporal prior knowledge can enable the model to understand dynamic semantics. Thus, our Temporal knowledge Compensator (TKC) aggregates essential semantic information by injecting temporal prior knowledge into visual features.

Specifically, we obtain global semantic vector $p_{g} \in \mathbb{R}^{1 \times C}$ by pooling temporal attribute features, and add it to frame-level features $[f_1, f_2, ..., f_T] \in \mathbb{R}^{T \times C}$:

$$\boldsymbol{F}_{\mathrm{q}} = [\boldsymbol{f}_{1} + \boldsymbol{p}_{\mathrm{g}}, ..., \boldsymbol{f}_{T} + \boldsymbol{p}_{\mathrm{g}}], \qquad (3)$$

where $F_q \in \mathbb{R}^{T \times C}$ is the obtained frame-level prototypes, which incorporate overall semantic information. The frame-level prototypes further aggregate temporal context information from temporal prior knowledge via vision and attribute cross-attention mechanism. Then the frame prototypes are fed into the temporal transformer [16] to capture the temporal relationships between frame-level prototypes. This is given by

$$\boldsymbol{P}_{t} = \mathtt{Tformer}(\mathtt{Softmax}(\boldsymbol{F}_{q}\boldsymbol{K}_{t}^{\top})\boldsymbol{V}_{t} + \boldsymbol{F}_{q}), \tag{4}$$

where K_t and V_t are the linear transformation features of temporal attribute features Q_t , Tformer is the temporal transformer [16], $P_t \in \mathbb{R}^{T \times C}$ is the frame-level prototypes capturing action dynamic information. In this way, the learned frame-level prototypes can adaptively perceive temporal changes and encode the action temporal context with the guidance of temporal knowledge.

3.6 Few-shot Metric

278

279 280

281

282 283

284 285

286

287

288

292

298

299 300

301

302 303

304

311 312 313

318

319

Few-shot Spatial Metric. To conduct spatial feature matching between videos, we propose an object-level prototype matching strategy based on the bidirectional Hausdorff Distance [16], which calculates the distances between query object-level prototypes and support object-level prototypes from the set matching perspective. Specifically, given query object-level prototypes $P_s \in \mathbb{R}^{T \times N \times C}$ and support object-level prototypes $\hat{P}_s \in \mathbb{R}^{T \times N \times C}$, we apply a bidirectional Mean Hausdorff metric to obtain a frame-level distance matrix $\hat{D} = [d_{ij}]_{T \times T} \in \mathbb{R}^{T \times T}$ as:

$$d_{ij} = \frac{1}{N} \sum_{\boldsymbol{p}_{i,k}^{s} \in \boldsymbol{p}_{i}^{s}} (\min_{\hat{\boldsymbol{p}}_{j,l}^{s} \in \hat{\boldsymbol{p}}_{j}^{s}} \| \boldsymbol{p}_{i,k}^{s} - \hat{\boldsymbol{p}}_{j,l}^{s} \|) + \frac{1}{N} \sum_{\hat{\boldsymbol{p}}_{j,l}^{s} \in \hat{\boldsymbol{p}}_{j}^{s}} (\min_{\boldsymbol{p}_{i,k}^{s} \in \boldsymbol{p}_{i}^{s}} \| \hat{\boldsymbol{p}}_{j,l}^{s} - \boldsymbol{p}_{i,k}^{s} \|),$$
(5)

where $p_{i,j}^s$ and $\hat{p}_{i,j}^s$ are the *j*-th support and query object-level prototypes in *i*-th frame, respectively. Then, for the frame-level distance matrix $\hat{D} \in \mathbb{R}^{T \times T}$, we find the smallest distance across the frame sequences, which gives a more confident probability of spatial feature matching. Finally, the spatial metric can be formulated as:

$$\mathcal{D}_{s} = \frac{1}{T} \sum_{i=1}^{T} (\min_{d_{i,j} \in \hat{\mathcal{D}}} \| d_{ij} \|) + \frac{1}{T} \sum_{j=1}^{T} (\min_{d_{i,j} \in \hat{\mathcal{D}}} \| d_{ij} \|).$$
(6)

Few-shot Temporal Metric. After obtaining the frame-level prototypes of support and query videos in a few-shot task, like in previous works [11, 24], we obtain support-query matching results by applying the temporal alignment metric:

$$\mathcal{D}_t = \texttt{Metric}(\boldsymbol{P}_t, \hat{\boldsymbol{P}}_t),$$
 (7)

325	settings are conducted under	nucled under the 5-way X-shot. Inter-Kiyoo denotes Keshet-50 pre-trained on imagenet.								
206	Mathad	Dro training	Bro training HMDB51				UCF101			
320	Wethou	i ic-training	1-shot	3-shot	5-shot	1-shot	3-shot	5-shot		
327	ARN [10] [ECCV20]	C3D	45.5	-	60.6	66.3	-	83.1		
328	OTAM [11] [CVPR20]	INet-RN50	54.5	65.7	-	79.9	87.0	-		
000	TRX [18] [CVPR21]	INet-RN50	53.1	66.8	75.6	78.2	92.4	96.1		
329	MTFAN [23] [CVPR22]	INet-RN50	59.0	-	74.6	84.8	-	95.1		
330	HyRSM [16] [CVPR22]	INet-RN50	60.3	71.7	76.0	83.9	93.0	94.7		
221	STRM [15] [CVPR22]	INet-RN50	52.3	67.4	77.3	80.5	92.7	96.9		
331	CPM [17] [ECCV22]	INet-RN50	60.1	-	-	71.4	-	-		
332	HCL [19] [ECCV22]	INet-RN50	59.1	71.2	76.3	82.6	91.0	94.5		
333	MoLo [70] [CVPR23]	INet-RN50	60.8	72.0	77.4	86.0	93.5	95.5		
555	GgHM [13] [ICCV23]	INet-RN50	61.2	-	76.9	85.2	-	96.3		
334	CLIP-FSAR [24] [IJCV24]	CLIP-RN50	69.2	77.6	80.3	91.3	95.1	97.0		
335	CapFSAR [21] [Arxiv23]	BLIP-ViT-B	65.2	-	78.6	93.3	-	97.8		
000	CLIP-Freeze [27] [ICML21]	CLIP-ViT-B	58.2	72.7	77.0	89.7	94.3	95.7		
336	CLIP-FSAR [24] [IJCV24]	CLIP-ViT-B	75.8	84.1	87.7	96.6	98.4	99.0		
337	DIST (Ours)	CLIP-ViT-B	82.6±0.3	87.1 ± 0.3	88.7 ± 0.1	98.3±0.2	99.0±0.2	99.2 ±0.1		

Table 1: Quantitative comparison results on HMDB51 [31] and UCF101 [32] (see §4.3). The experiment settings are conducted under the 5-way *K*-shot. "INet-RN50" denotes ResNet-50 pre-trained on ImageNet.

where $P_t \in \mathbb{R}^{T \times C}$ represents query frame-level prototypes, $\hat{P}_t \in \mathbb{R}^{T \times C}$ is support frame-level prototypes, and Metric denotes the OTAM [11] metric by default. We formulate the distance between support and query videos as the weighted sum of the distances obtained by the few-shot spatial metric and few-shot temporal metric:

$$\mathcal{D} = \mathcal{D}_{\rm t} + \alpha \mathcal{D}_{\rm s},\tag{8}$$

where α is a coefficient parameter. Our proposed matching strategy combines the advantages of frame- and object-level prototype matching to cope with appearance- and motion-centric actions.

Following previous works [11, 18, 16], we minimize cross-entropy loss \mathcal{L}_{CE} over the support-query distances based on the ground-truth labels to end-to-end train DIST. For few-shot inference, total support-query distance in Eq. 8 is employed as logits to produce final query prediction.

4 EXPERIMENTS

338

343

347

348

349 350 351

352 353

354

361

362

363 364

4.1 EXPERIMENTAL SETUP

Dataset. We conduct extensive experiments on five datasets, *i.e.*, Kinetics [33], SSv2-full [34],
SSv2-small [34], HMDB51 [31], and UCF101 [32]. For SSv2-full [34], SSv2-Small [34] and Kinetics [33], we utilize the split as in CMN [52], with 64, 12, and 24 classes used for train, val,
and test, respectively. For HMDB51 [31] and UCF101 [32], we adopt the split setting as in
ARN [10], where the 51 classes in HMDB51 are split into 31/10/10 classes for train/val/test,
while the 101 classes in UCF101 are split into 70/10/21 classes for train/val/test.

Evaluation. Following the official evaluation protocols [13, 11], we use 5-way 1-shot and 5-shot accuracy for evaluation, and report average results over 10,000 tasks randomly selected from test.

4.2 IMPLEMENTATION DETAILS

Network Architecture. We use CLIP ViT-B [27] as our backbone for a fair comparison with previous methods [24, 21]. By default, the number of spatial attributes G and temporal attributes Lare set to 6 and 3, respectively (ablation study in Table 7 of Appendix). The number of object-level prototypes N is 9. The value of parameter α is set to 0.5 (see Fig. 4 (left)).

Network Training. Following previous methods [1, 16, 11, 13], we uniformly and sparsely sample T = 8 frames of each video to encode video representation. In the training phase, we adopt basic data augmentation, such as random horizontal flipping, cropping, and color jitter. To retain the original pre-trained prior knowledge in the text encoder and reduce the optimization burden, we freeze the text encoder and prevent it from being updated during training. Moreover, we use the Adam [71] optimizer with the multi-step scheduler to train our framework.

Reproducibility. DIST is implemented in PyTorch, and all models are trained and tested on two NVIDIA Tesla V100 GPUs with a 32GB memory per card. Full code will be released.

380	Mathad	Des teriaines	Kine	etics	SS	v2	SSv2-	small
381	Method	Pre-training	1-shot	5-shot	1-shot	5-shot	1-shot	5-shot
200	ARN [10] [ECCV20]	C3D	63.7	82.4	-	-	-	-
302	OTAM [11] [CVPR20]	INet-RN50	73.0	85.8	42.8	52.3	36.4	48.0
383	TRX [18] [CVPR21]	INet-RN50	63.6	85.9	42.0	64.6	36.0	56.7
38/	MTFAN [23] [CVPR22]	INet-RN50	74.6	87.4	45.7	60.4	-	-
304	HyRSM [16] [CVPR22]	INet-RN50	73.7	86.1	54.3	69.0	40.6	56.1
385	STRM [15] [CVPR22]	INet-RN50	62.9	86.7	43.1	68.1	37.1	55.3
386	CPM [17] [ECCV22]	INet-RN50	73.3	-	49.3	66.7	-	-
000	HCL [19] [ECCV22]	INet-RN50	73.7	85.8	47.3	64.9	38.9	55.4
387	MoLo [70] [CVPR23]	INet-RN50	74.0	85.6	56.6	70.6	42.7	56.4
388	CLIP-FSAR [24] [IJCV24]	CLIP-RN50	87.6	91.9	58.1	62.8	52.0	55.8
000	CapFSAR [21] [Arxiv23]	BLIP-ViT-B	84.9	93.1	51.9	68.2	45.9	59.9
389	CLIP-Freeze [27] [ICML21]	CLIP-ViT-B	78.9	91.9	30.0	42.4	29.5	42.5
390	CLIP-FSAR [24] [IJCV24]	CLIP-ViT-B	89.7	95.0	61.9	72.1	54.5	61.8
391	DIST (Ours)	CLIP-ViT-B	92.7 ± 0.3	95.5 ± 0.1	64.2±0.2	75.2 ± 0.2	57.5 ± 0.3	62.5 ± 0.1

378 Table 2: Quantitative comparison results on Kinetics [33] and SSv2-small [34] (see §4.3). The experiment settings are conducted under the 5-way K-shot. "INet-RN50" denotes ResNet-50 pre-trained on ImageNet. 379

Table 3: Comparison results [34, 32, 31, 33] (see §4.3) by combining few-shot and zero-shot results.

Madaad	HMI	DB51	UCI	F101	Kin	etics	SS	v2	SSv2-	-small	
	Method	1-shot	5-shot	1-shot	5-shot	1-shot	5-shot	1-shot	5-shot	1-shot	5-shot
	CLIP-FSAR [24] [IJCV24]	77.1	87.7	97.0	99.1	94.8	95.4	62.1	72.1	54.6	61.8
	DIST (Ours)	82.6	88.7	98.3	99.2	95.6	96.0	64.6	75.8	57.5	62.5

³⁹⁶ 397 398

39 392

393

4.3 COMPARISON WITH STATE-OF-THE-ARTS

399 We compare the performance of our DIST with current state-of-the-art FSAR methods on five stan-400 dard datasets [31, 32, 34, 33] in Table 1 and Table 2. It demonstrates that DIST outperforms all 401 FSAR methods. Specifically, compared to CLIP-FSAR [24] that only uses naive class names as 402 semantic information, our approach achieves better results in multiple datasets and task settings. It indicates our DIST further boosts performance by grasping spatiotemporal-decoupled prior knowl-403 edge from LLM to compensate for visual features. Further, the performance margin between DIST 404 and CLIP-FSAR is more significant under low shots. Notably, on HMDB51 [31] and UCF101 [32] 405 datasets, the performance of our DIST on the 5-way 3-shot setting is comparable to the performance 406 of CLIP-FSAR on the 5-way 5-shot setting. We attribute this to the fact that decoupled prompts 407 provide enough background knowledge to compensate for visual features than naive class names es-408 pecially when visual information is insufficient (*i.e.*, one shot). Furthermore, Table 3 compares our 409 DIST against CLIP-FSAR under another setting, which makes few-shot predictions with the help 410 of zero-shot results. The results show DIST consistently outperforms CLIP-FSAR on each dataset.

411 412 413

4.4 ABLATION STUDY

414 We conduct ablation experiments to evaluate the efficacy of our idea and core model designs. Unless 415 otherwise specified, we adopt CLIP-ViT-B model as default experimental setting.

416 Key Component Analysis. Table 4 summa-417 rizes the impact of each module in DIST. 418 Specifically, compared to the baseline, Tempo-419 ral Knowledge Compensator (TKC) (cf.§3.5) 420 brings 5.2%, 2.3% and 1.9% performance 421 gains on HMDB51 [31], SSv2-small [33] and 422 UCF101 [32], respectively. This consistent promotion indicates that TKC can enhance the tem-423

Table	4:	Impacts	of	core	coi	nponents	s on
HMDB	51 [<mark>31</mark>]	, SSv2-sr	nall	[33],	and	UCF101	[32]
in the 5-	-way 1-	shot tasks	(see	§4.4)			

•			
Method Component	HMDB51	SSv2-small	UCF101
BASELINE	75.8	53.8	96.0
SKC only	77.6	55.7	96.6
TKC only	81.0	56.1	97.9
DIST (Ours)	82.6	57.5	98.3

poral awareness of DIST to facilitate accurate matching. In addition, the proposed Spatial Knowl-424 edge Compensator (SKC) (cf.§3.4) improves on the three datasets [31, 32] by 1.8%, 1.9% and 425 0.6%, respectively, which indicates leveraging spatial prior knowledge can focus on action-related 426 spatial details to boost few-shot performance. Moreover, combining the two modules can further im-427 prove performance, indicating the complementarity between spatial and temporal prior knowledge. 428

429 Attribute Injection Manners. We respectively propose SKC (cf.§3.4) and TKC (cf.§3.5) to inject spatial/temporal attribute features into visual features. In Table 5a, we study the effect of different 430 temporal/spatial attribute injection manners. "Concat" means that the visual features and attribute 431 features are directly concatenated and then fed into transformers for multimodal fusion like CLIP-

433	Spatia	A think	Tamp	anal Attribute		DD51	Cna	tial Attail	auto	Tamana	nal Attaibute	LIME	DD51
40.4	Comoo	t SVC	e Tempo	TUC		5 abot	I spai	Va and	Jule Ladra	I ab al	Va anda daa		JDJ1 Fahat
434	Conca	u src	Conca		1-shot	o-snot	Laber	Know	leage	Laber	Knowledge	1-snot	5-snot
435	1		✓		80.7	88.3	1			1		80.0	87.3
400		1	✓		81.0	88.6		~	·	1		81.2	88.0
430	1			1	81.6	88.5	1				1	81.6	88.6
437					82.6	88.7					 Image: A second s	82.6	88.7
438		(a)	attribute	injection ma	nner				(b)	attribu	te content		
439	Me	ethod	HMDB51	SSv2-small	Carefal M	- 4 ¹	HMI	DB51	Tem	ooral Me	tric	HM	DB51
440		0 4 D 70 (1		7 0.0	Spatial M	etric	1-shot	5-shot	Temp		une	1-shot	5-shot
	CLIP-F	SAR [24]	75.8	53.8	One-to-on	e matching	82.4	87.8	CLIF	P-FSAR (Bi-MHM) [24] 76.0	87.8
441	CLIP-	-FSAR [†]	81.0	56.1	Bi-MHM	[16]	82.4	87.9	Ours	S (Bi-ME	IM)	82.7	88.9
442	DIST	(Ours)	82.6	57.5	Ours	[10]	82.6	88.7	Ours	-F5AK (ε (ΟΤΔΝ	OIAM) [24] D	75.8 82.6	81.1 88.7
1/13	(c)	cnowled	ge comp	ensator	(1)	. 1 . 1					1 4 1 *		
	(0)1	anowieu	se comp	chisator	(d) spa	liai match	ing me	etric	((e) temp	boral matchi	ng metr	10
444													
445	83.0 _L					100							
4.4.0	02.5					100				_		- (CLIP-FSAR
440	82.3				ି	90	-						Durs
447	\$ 82.0			•	e),	80							_
1/18	lac				ac	70							_
0	∃ 81.5	/			cur	60							
449	V al a				Ac	00			_				
450	· 81.0					50							
454	80.5					40 💻		11	N				11
451	2010	0 0.1	0.3 0.5	5 0.7 1		renci	ng kick	ick ball	pick	pour nist	nh un en	smoke	talk
450		The fi	usion para	umeter α		101		KIC		· pu		<u>.</u>	

Table 5: A set of ablation studies on HMDB51 [31] (see §4.4). The adopted network designs are marked in red.

Figure 4: Left: The impact of the varying fusion parameter α on HMDB51 [31] in the 5-way 1-shot setting (see §4.4). Right: 5-way 1-shot class improvement of DIST compared to CLIP-FSAR [24] on all class action classes on HMDB51 [31] (see §4.5). Our DIST achieves improvement on all action classes.

FSAR [24]. The experimental results show that our proposed SKC and TKC yields better results, suggesting the effectiveness of our module design.

Attribute Content. We investigate the impact of different temporal and spatial attribute content on the performance of our proposed DIST. As shown in Table 5b, we observe that utilizing spatial and temporal prior knowledge generated by LLM consistently performs better than using class names, with 1.2% and 1.6% performance gains in the 1-shot setting on HMDB [31], respectively. In addition, combining spatial and temporal prior knowledge yields better results, which demonstrate different prior knowledge is complementary to others.

Impact of Knowledge Compensators. We replace the category labels of CLIP-FSAR [24] with LLM-generated prompts (*i.e.*, CLIP-FSAR^{\dagger}) and report the comparison results in Table 5c. DIST gains larger improvements Compared to CLIP-FSAR[†]. This suggests our performance gains are not solely due to the usage of LLM-generated prompts, but also due to proposed knowledge compen-sators which make full use of LLM prompts to compensate for visual features.

Matching Metrics. We analyze the impact of different spatial matching metrics in Table 5d. We adopt different spatial matching metrics (cf. Eq. 6), including one-to-one matching, Bi-MHM [16], and our proposed spatial metric. One-to-one matching means computing the spatial matching scores of aligned object-level prototypes between the support video and query video. The results show that our proposed spatial metric achieves the best results, suggesting the effectiveness of our proposed spatial matching metric. We also conduct experiments using different temporal matching metrics (cf. Eq. 7) on HMDB51 [31]. As shown in Table 5e, our method can adapt to any temporal alignment metric and achieves better performance compared to CLIP-FSAR [24].

Varying Fusion Parameter α . Fig. 4 (left) shows the impact of the varying fusion parameter α (cf. Eq. 8) of spatial and temporal matching in the 5-way 1-shot task on HMDB51 [31]. From the results, the optimal value of parameter α is 0.5 for HMDB51.

4.5 QUALITY ANALYSIS

Class-wise Performance Gains. Fig. 4 (right) shows 5-way 1-shot class-wise performance gains obtained by our DIST over CLIP-FSAR [24] on HMDB51 [31]. Notably, our DIST achieves perfor-mance gains in all action classes. We also observe that DIST achieves gains above 10% for classes



Figure 6: Visualization of spatial and temporal prompts under the 5-way 1-shot setting (see §4.5). The spatial prompts are shown as highlighted response areas in each frame. We also show cross-attention temporal prompt weights of Eq. 4 in a line graph.

such as *run*, *pour*, *kick ball*, *etc*.. It indicates that the spatiotemporal-decoupled prior knowledge can
 easily include objects involved in these actions and capture action-related dynamic information.

Visualization of Feature Distribution. To further qualitatively analyze the changes in feature distribution after incorporating spatiotemporal-decoupled prior knowledge, we follow previous methods [24, 15] to visualize the feature distribution of CLIP-FSAR [24] (only category name as semantic information) and our algorithm DIST in Fig. 5. We observe that after utilizing action-related prior knowledge, our method shows more compact intra-class feature distributions and more discriminative inter-class features. As shown in Fig. 5 (a), the three classes "run", "fencing", and "kick" become clearly distinguishable from each other after injecting decoupled attributes into visual features.

520 Visualization of Spatial and Temporal attributes. To analyze the role of spatial and temporal 521 attributes in DIST, Fig. 6 displays the visualization results of these attributes. As seen, the attention maps of our DIST focus more on action-related objects and reduce attention to the background 522 and unrelated objects. This demonstrates our DIST grasps prior knowledge provided by spatial 523 attributes to capture spatial details. Then, we calculate the cross-attention scores between temporal 524 attributes and frames according to Eq. 4. It can be seen that different temporal attributes have 525 different weights on the frame sequences, which proves that our DIST can learn temporal relations 526 and capture dynamic semantics. For example, the temporal attribute "Hold container" has larger 527 weights on the first three frames, which indicates these frames may correspond to dynamic semantics 528 implied by the temporal attribute. See more examples in C_2 of Appdendix.

529 530 531

532

509

510

5 CONCLUSION

In this work, we propose a novel yet effective DIST framework for FSAR, which is the first work to grasp spatiotemporal-decoupled prior knowledge from LLM to compensate for visual features.
In particular, we design Spatial/Temporal Knowledge Compensators to learn object- and frame-level prototypes, so as to capture fine-grained spatial details and dynamic semantics. Experimental results demonstrate that our DIST achieves state-of-the-art performance on four standard bench-marks. However, this is the first cursory exploration of leveraging spatiotemporal-decoupled prior knowledge in FSAR. Though the first step is not always elegant, exploring additional attempts to leverage richer prior knowledge provided by LLMs promises intriguing prospects for the future.

540 REFERENCES

547

548

549 550

551

552

553

554

555 556

558

559

560 561 562

563

564 565

566

567

568 569

570

571 572

573

574

575

576

577

578 579

580

581 582

583

584

585

586

587

588 589

590

591

- [1] Limin Wang, Yuanjun Xiong, Zhe Wang, Yu Qiao, Dahua Lin, Xiaoou Tang, and Luc Van Gool. Temporal segment networks: Towards good practices for deep action recognition. In *ECCV*, pages 20–36, 2016. 1, 7
- Zitian Chen, Yanwei Fu, Yu-Xiong Wang, Lin Ma, Wei Liu, and Martial Hebert. Image deformation meta-networks for one-shot learning. In *CVPR*, pages 8680–8689, 2019. 1, 3
 - [3] Ji Lin, Chuang Gan, and Song Han. Tsm: Temporal shift module for efficient video understanding. In *ICCV*, pages 7083–7093, 2019.
 - [4] Limin Wang, Bingkun Huang, Zhiyu Zhao, Zhan Tong, Yinan He, Yi Wang, Yali Wang, and Yu Qiao. Videomae v2: Scaling video masked autoencoders with dual masking. In *CVPR*, pages 14549–14560, 2023. 1
 - [5] Ceyuan Yang, Yinghao Xu, Jianping Shi, Bo Dai, and Bolei Zhou. Temporal pyramid network for action recognition. In *CVPR*, pages 591–600, 2020. 1
 - [6] Limin Wang, Zhan Tong, Bin Ji, and Gangshan Wu. Tdn: Temporal difference networks for efficient action recognition. In *CVPR*, pages 1895–1904, 2021. 1
 - [7] Zhen Xing, Qi Dai, Han Hu, Jingjing Chen, Zuxuan Wu, and Yu-Gang Jiang. Svformer: Semisupervised video transformer for action recognition. In CVPR, pages 18816–18826, 2023.
 - [8] Yifei Huang, Lijin Yang, Guo Chen, Hongjie Zhang, Feng Lu, and Yoichi Sato. Matching compound prototypes for few-shot action recognition. *International Journal of Computer Vision*, pages 1–26, 2024. 1
 - [9] Jintao Lin, Zhaoyang Liu, Wenhai Wang, Wayne Wu, and Limin Wang. Vlg: General video recognition with web textual knowledge. *International Journal of Computer Vision*, pages 1–26, 2024. 1
 - [10] Hongguang Zhang, Li Zhang, Xiaojuan Qi, Hongdong Li, Philip HS Torr, and Piotr Koniusz. Few-shot action recognition with permutation-invariant attention. In *ECCV*, pages 525–542, 2020. 1, 7, 8
 - [11] Kaidi Cao, Jingwei Ji, Zhangjie Cao, Chien-Yi Chang, and Juan Carlos Niebles. Few-shot video classification via temporal alignment. In *CVPR*, pages 10618–10627, 2020. 1, 3, 6, 7, 8, 16
 - [12] Shuyuan Li, Huabin Liu, Rui Qian, Yuxi Li, John See, Mengjuan Fei, Xiaoyuan Yu, and Weiyao Lin. Ta2n: Two-stage action alignment network for few-shot action recognition. In AAAI, volume 36, pages 1404–1411, 2022. 1
 - [13] Jiazheng Xing, Mengmeng Wang, Yudi Ruan, Bofan Chen, Yaowei Guo, Boyu Mu, Guang Dai, Jingdong Wang, and Yong Liu. Boosting few-shot action recognition with graph-guided hybrid matching. In *ICCV*, pages 1740–1750, 2023. 1, 3, 4, 7, 16, 19
 - [14] Haifeng Xia, Kai Li, Martin Renqiang Min, and Zhengming Ding. Few-shot video classification via representation fusion and promotion learning. In *ICCV*, pages 19311–19320, 2023.
 - [15] Anirudh Thatipelli, Sanath Narayan, Salman Khan, Rao Muhammad Anwer, Fahad Shahbaz Khan, and Bernard Ghanem. Spatio-temporal relation modeling for few-shot action recognition. In *CVPR*, pages 19958–19967, 2022. 1, 3, 5, 7, 8, 10, 19
 - [16] Xiang Wang, Shiwei Zhang, Zhiwu Qing, Mingqian Tang, Zhengrong Zuo, Changxin Gao, Rong Jin, and Nong Sang. Hybrid relation guided set matching for few-shot action recognition. In *CVPR*, pages 19948–19957, 2022. 1, 3, 4, 6, 7, 8, 9, 16, 19
- 593 [17] Yifei Huang, Lijin Yang, and Yoichi Sato. Compound prototype matching for few-shot action recognition. In ECCV, pages 351–368, 2022. 1, 7, 8

600

601

602

603

604

605

606

607

608 609

610

611

616

617

618 619

620

621

622

623

624

625

626

627

628

629 630

631

632 633

634

635

636

637

638

639 640

641

642

643

- [18] Toby Perrett, Alessandro Masullo, Tilo Burghardt, Majid Mirmehdi, and Dima Damen. Temporal-relational crosstransformers for few-shot action recognition. In *CVPR*, pages 475–484, 2021. 1, 3, 7, 8
 - [19] Sipeng Zheng, Shizhe Chen, and Qin Jin. Few-shot action recognition with hierarchical matching and contrastive learning. In *ECCV*, pages 297–313, 2022. 1, 7, 8, 16
 - [20] Khoi D Nguyen, Quoc-Huy Tran, Khoi Nguyen, Binh-Son Hua, and Rang Nguyen. Inductive and transductive few-shot video classification via appearance and temporal alignments. In ECCV, pages 471–487, 2022. 1
 - [21] Xiang Wang, Shiwei Zhang, Hangjie Yuan, Yingya Zhang, Changxin Gao, Deli Zhao, and Nong Sang. Few-shot action recognition with captioning foundation models. arXiv preprint arXiv:2310.10125, 2023. 1, 3, 7, 8
 - [22] Oriol Vinyals, Charles Blundell, Timothy Lillicrap, Daan Wierstra, et al. Matching networks for one shot learning. In *NeurIPS*, volume 29, 2016. 1, 3
 - [23] Jiamin Wu, Tianzhu Zhang, Zhe Zhang, Feng Wu, and Yongdong Zhang. Motion-modulated temporal fragment alignment network for few-shot action recognition. In CVPR, pages 9151– 9160, 2022. 1, 3, 7, 8
- [24] Xiang Wang, Shiwei Zhang, Jun Cen, Changxin Gao, Yingya Zhang, Deli Zhao, and Nong Sang. Clip-guided prototype modulating for few-shot action recognition. *International Journal of Computer Vision*, 2023. 1, 3, 5, 6, 7, 8, 9, 10, 16, 18
 - [25] Hongyu Qu, Rui Yan, Xiangbo Shu, Haoliang Gao, Peng Huang, and Guo-Sen Xie. Mvpshot: Multi-velocity progressive-alignment framework for few-shot action recognition. arXiv preprint arXiv:2405.02077, 2024. 1
 - [26] Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc Le, Yun-Hsuan Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning with noisy text supervision. In *ICML*, pages 4904–4916, 2021. 1
 - [27] 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*, pages 8748–8763, 2021. 1, 3, 4, 7, 8
 - [28] Bolin Ni, Houwen Peng, Minghao Chen, Songyang Zhang, Gaofeng Meng, Jianlong Fu, Shiming Xiang, and Haibin Ling. Expanding language-image pretrained models for general video recognition. In ECCV, pages 1–18, 2022. 1
 - [29] Mengmeng Wang, Jiazheng Xing, Jianbiao Mei, Yong Liu, and Yunliang Jiang. Actionclip: Adapting language-image pretrained models for video action recognition. *IEEE Transactions* on Neural Networks and Learning Systems, 2023. 1
 - [30] Wenhao Wu, Zhun Sun, Yuxin Song, Jingdong Wang, and Wanli Ouyang. Transferring visionlanguage models for visual recognition: A classifier perspective. *International Journal of Computer Vision*, 132(2):392–409, 2024. 1
 - [31] Hildegard Kuehne, Hueihan Jhuang, Estíbaliz Garrote, Tomaso Poggio, and Thomas Serre. Hmdb: a large video database for human motion recognition. In *ICCV*, pages 2556–2563, 2011. 3, 7, 8, 9, 16, 17, 19
 - [32] 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. **3**, **7**, **8**, **16**, **17**
 - [33] Joao Carreira and Andrew Zisserman. Quo vadis, action recognition? a new model and the kinetics dataset. In *CVPR*, pages 6299–6308, 2017. 3, 7, 8, 16, 17, 19
- [34] Raghav Goyal, Samira Ebrahimi Kahou, Vincent Michalski, Joanna Materzynska, Susanne
 Westphal, Heuna Kim, Valentin Haenel, Ingo Fruend, Peter Yianilos, Moritz Mueller-Freitag,
 et al. The" something something" video database for learning and evaluating visual common sense. In *ICCV*, pages 5842–5850, 2017. 3, 7, 8, 16, 17

652

653 654

655

656

657

658

659 660

661

662 663

664

665

666

667

668

669 670

671

672

673

674 675

676

677

678

679

680

681

682 683

684

685

686

687 688

689

690

691

692

693 694

696 697

698

- [35] Li Fei-Fei, Robert Fergus, and Pietro Perona. One-shot learning of object categories. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 28(4):594–611, 2006. 3
 - [36] Yunhui Guo, Noel C Codella, Leonid Karlinsky, James V Codella, John R Smith, Kate Saenko, Tajana Rosing, and Rogerio Feris. A broader study of cross-domain few-shot learning. In ECCV, pages 124–141, 2020. 3
 - [37] Kai Li, Yulun Zhang, Kunpeng Li, and Yun Fu. Adversarial feature hallucination networks for few-shot learning. In *CVPR*, pages 13470–13479, 2020. 3
 - [38] Han-Jia Ye, Hexiang Hu, and De-Chuan Zhan. Learning adaptive classifiers synthesis for generalized few-shot learning. *International Journal of Computer Vision*, 129:1930–1953, 2021. 3
 - [39] Alexander J Ratner, Henry Ehrenberg, Zeshan Hussain, Jared Dunnmon, and Christopher Ré. Learning to compose domain-specific transformations for data augmentation. In *NeurIPS*, volume 30, 2017. 3
 - [40] Zitian Chen, Yanwei Fu, Yinda Zhang, Yu-Gang Jiang, Xiangyang Xue, and Leonid Sigal. Semantic feature augmentation in few-shot learning. arXiv preprint arXiv:1804.05298, 86(89):2, 2018. 3
 - [41] Zitian Chen, Yanwei Fu, Yinda Zhang, Yu-Gang Jiang, Xiangyang Xue, and Leonid Sigal. Multi-level semantic feature augmentation for one-shot learning. *IEEE Transactions on Image Processing*, 28(9):4594–4605, 2019. 3
 - [42] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *ICML*, pages 1126–1135, 2017. 3
 - [43] Muhammad Abdullah Jamal and Guo-Jun Qi. Task agnostic meta-learning for few-shot learning. In CVPR, pages 11719–11727, 2019. 3
 - [44] Aravind Rajeswaran, Chelsea Finn, Sham M Kakade, and Sergey Levine. Meta-learning with implicit gradients. In *NeurIPS*, volume 32, 2019. 3
 - [45] Andrei A Rusu, Dushyant Rao, Jakub Sygnowski, Oriol Vinyals, Razvan Pascanu, Simon Osindero, and Raia Hadsell. Meta-learning with latent embedding optimization. arXiv preprint arXiv:1807.05960, 2018. 3
 - [46] Sung Whan Yoon, Jun Seo, and Jaekyun Moon. Tapnet: Neural network augmented with task-adaptive projection for few-shot learning. In *ICML*, pages 7115–7123. PMLR, 2019. 3
 - [47] Jake Snell, Kevin Swersky, and Richard Zemel. Prototypical networks for few-shot learning. In *NeurIPS*, volume 30, 2017. 3
 - [48] Han-Jia Ye, Hexiang Hu, De-Chuan Zhan, and Fei Sha. Few-shot learning via embedding adaptation with set-to-set functions. In *CVPR*, pages 8808–8817, 2020. 3
 - [49] Flood Sung, Yongxin Yang, Li Zhang, Tao Xiang, Philip HS Torr, and Timothy M Hospedales. Learning to compare: Relation network for few-shot learning. In CVPR, pages 1199–1208, 2018. 3
 - [50] Hongyang Li, David Eigen, Samuel Dodge, Matthew Zeiler, and Xiaogang Wang. Finding task-relevant features for few-shot learning by category traversal. In *CVPR*, pages 1–10, 2019.
 3
 - [51] Yizhou Zhao, Zhenyang Li, Xun Guo, and Yan Lu. Alignment-guided temporal attention for video action recognition. In *NeurIPS*, volume 35, pages 13627–13639, 2022. 3
 - [52] Linchao Zhu and Yi Yang. Compound memory networks for few-shot video classification. In *ECCV*, pages 751–766, 2018. **3**, **7**, **16**
- [53] Ning Ma, Hongyi Zhang, Xuhui Li, Sheng Zhou, Zhen Zhang, Jun Wen, Haifeng Li, Jingjun Gu, and Jiajun Bu. Learning spatial-preserved skeleton representations for few-shot action recognition. In *ECCV*, pages 174–191. Springer, 2022. 3

713 714

715

716

717

718

719

720 721

722

723 724

725 726

727

728 729

730

731

732

733

734 735

736

737

738

739 740

741

742

743

744

745 746

747

- [54] Yuqian Fu, Li Zhang, Junke Wang, Yanwei Fu, and Yu-Gang Jiang. Depth guided adaptive meta-fusion network for few-shot video recognition. In *ACM MM*, pages 1142–1151, 2020. 3
- [55] Linchao Zhu and Yi Yang. Label independent memory for semi-supervised few-shot video classification. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(1):273–285, 2020. 3
- [56] Xiao Wang, Weirong Ye, Zhongang Qi, Xun Zhao, Guangge Wang, Ying Shan, and Hanzi Wang. Semantic-guided relation propagation network for few-shot action recognition. In *ACM MM*, pages 816–825, 2021. 3
 - [57] Jiazheng Xing, Mengmeng Wang, Yong Liu, and Boyu Mu. Revisiting the spatial and temporal modeling for few-shot action recognition. In AAAI, volume 37, pages 3001–3009, 2023. 3, 5
 - [58] Xiao Wang, Weirong Ye, Zhongang Qi, Guangge Wang, Jianping Wu, Ying Shan, Xiaohu Qie, and Hanzi Wang. Task-aware dual-representation network for few-shot action recognition. IEEE Transactions on Circuits and Systems for Video Technology, 33(10):5932–5946, 2023. 3
 - [59] Jiazheng Xing, Mengmeng Wang, Xiaojun Hou, Guang Dai, Jingdong Wang, and Yong Liu. Multimodal adaptation of clip for few-shot action recognition. arXiv preprint arXiv:2308.01532, 2023. 3, 5
 - [60] Zhimao Peng, Zechao Li, Junge Zhang, Yan Li, Guo-Jun Qi, and Jinhui Tang. Few-shot image recognition with knowledge transfer. In *ICCV*, pages 441–449, 2019. 3
 - [61] Jingyi Xu and Hieu Le. Generating representative samples for few-shot classification. In *CVPR*, pages 9003–9013, 2022. 3
 - [62] Wentao Chen, Chenyang Si, Zhang Zhang, Liang Wang, Zilei Wang, and Tieniu Tan. Semantic prompt for few-shot image recognition. In *CVPR*, pages 23581–23591, 2023. 3, 4
 - [63] Chen Xing, Negar Rostamzadeh, Boris Oreshkin, and Pedro O O Pinheiro. Adaptive crossmodal few-shot learning. In *NeurIPS*, volume 32, 2019. 3, 4
 - [64] Kun Yan, Chenbin Zhang, Jun Hou, Ping Wang, Zied Bouraoui, Shoaib Jameel, and Steven Schockaert. Inferring prototypes for multi-label few-shot image classification with word vector guided attention. In AAAI, volume 36, pages 2991–2999, 2022. 3
 - [65] Felix Wu, Amauri Souza, Tianyi Zhang, Christopher Fifty, Tao Yu, and Kilian Weinberger. Simplifying graph convolutional networks. In *ICML*, pages 6861–6871, 2019. 3
 - [66] Hai Zhang, Junzhe Xu, Shanlin Jiang, and Zhenan He. Simple semantic-aided few-shot learning. *arXiv preprint arXiv:2311.18649*, 2023. **3**, **4**
 - [67] Hao Tang, Jun Liu, Shuanglin Yan, Rui Yan, Zechao Li, and Jinhui Tang. M3net: Multi-view encoding, matching, and fusion for few-shot fine-grained action recognition. In ACM MM, pages 1719–1728, 2023. 4
 - [68] Tom B Brown. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*, 2020. 5
 - [69] OpenAI. Chatgpt, 2023. https://openai.com/blog/chatgpt/, Last accessed on 2024-01-13. 5
- [70] Xiang Wang, Shiwei Zhang, Zhiwu Qing, Changxin Gao, Yingya Zhang, Deli Zhao, and Nong Sang. Molo: Motion-augmented long-short contrastive learning for few-shot action recognition. In *CVPR*, pages 18011–18021, 2023. 7, 8
- [71] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980, 2014. 7, 16
- [72] Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of machine learning research*, 9(11), 2008. 10

[73] Ming Hu, Lin Wang, Siyuan Yan, Don Ma, Qingli Ren, Peng Xia, Wei Feng, Peibo Duan, Lie Ju, and Zongyuan Ge. Nurvid: A large expert-level video database for nursing procedure activity understanding. In NeurIPS, volume 36, 2024. 19 [74] Xiangbo Shu, Jiawen Yang, Rui Yan, and Yan Song. Expansion-squeeze-excitation fusion network for elderly activity recognition. IEEE Transactions on Circuits and Systems for Video Technology, 32(8):5281-5292, 2022. 19

This appendix provides additional details for the ICLR 2025 submission, titled "Spatio-temporal Decoupled Knowledge Compensator for Few-shot Action Recognition". The appendix is organized as follows:

- §A provides additional implementation details. 814 815 • §B provides the pseudo-code of spatial and temporal feature matching. 816 • §C introduces more quantitative and qualitative experiment results. 817 §D shows more additional examples of spatio-temporal knowledge. 818 • §E discusses our limitations and social impact. 819 820 821 Α ADDITIONAL IMPLEMENTATION DETAILS 822 823 A.1 DATASET DETAILS 824 We compare the proposed DIST with top-leading methods on four few-shot benchmarks, including 825 Kinetics [33], HMDB51 [31], UCF101 [32], and SSv2-small [34]. 826 827 • HMDB51 [31] contains 51 action classes and has 6,766 video clips. Following the setting 828 of previous methods [52, 24] for few-shot action recognition, we divide them into 31, 10, 829 and 10 action classes as train/val/test. The constructed dataset has 31 action classes with
 - for meta-training, meta-validation, and meta-testing.
 Kinetics [33] contains 400 action classes and has 306,245 video clips. Following the setting of previous methods [52, 16] for few-shot action recognition, we select 100 classes and divide them into 64, 12, and 24 action classes as train/val/test. The constructed dataset has 64 action classes with 6,389 videos, 12 action classes with 1,199 videos, and 24 action classes with 2,395 videos for meta-training, meta-validation, and meta-testing.

4,280 videos, 10 action classes with 1,194 videos, and 10 action classes with 1,292 videos

- UCF101 [32] contains 101 action classes and has 13,320 video clips. Following the setting of previous methods [52, 24] for few-shot action recognition, we divide them into 70, 10, and 21 action classes as train/val/test. The constructed dataset has 70 action classes with 9,154 videos, 10 action classes with 1,421 videos, and 10 action classes with 2,745 videos for meta-training, meta-validation, and meta-testing.
- **SSv2-small** [34] contains 174 action classes and has 220,847 video clips. Following the setting of previous methods [11, 24] for few-shot action recognition, we select 100 classes and divide them into 64, 12, and 24 action classes as the meta-training, meta-validation, and meta-testing set, respectively.
- A.2 IMPLEMENTATION DETAILS

848 Following previous methods [13, 24, 19, 16], we uniformly sample T = 8 frames from input videos, 849 which are scaled to a height of 256. In the training phase, we adopt basic data augmentation, such as 850 random horizontal flipping, cropping, and color jitter. In contrast, only a center crop is used during 851 the testing phase. We use the PyTorch library to train our DIST on two Tesla v100 GPUs. Moreover, 852 our framework uses the Adam optimizer [71] with the multi-step scheduler to train our model. The 853 total number of training steps is set to 10. Table 6 shows the same settings of hyperparameters with CLIP-FSAR [24] in a multi-step scheduler for various datasets. In this table, **Ir** represents the 854 learning rate, steps indicates the number of steps to change the learning rate, iter_st refers to the 855 number of iterations per step, and lr_st denotes the multiplication factor for updating the learning 856 rate at each changing step.

857 858

830

831 832

833

834

835

836

837 838

839

840

841

842

843

844

845 846

847

859 860

B PSEUDO CODE OF SPATIAL AND TEMPORAL MATCHING

Algorithm 1 provides the pseudo-code of spatial and temporal feature matching in a PyTorch-like
 style. With acquired action-related prior knowledge, we decompose video matching process into
 complementary object-level and frame-level prototype matching. To guarantee reproducibility, full
 code will be released.

872

899 900

901 902 903

904

905

906

907

908

909 910

Dataset	lr	steps	iter_st	lr_st
Kinetics [33]	1e - 5	[0,6,9]	1000	[1,0.1,0.01]
HMDB51 [31]	1e - 5	[0,4,6]	2800	[1,0.1,0.01]
UCF101 [32]	2e - 6	[0,4,6]	2400	[1,0.1,0.01]
SSv2-small [34]	5e-5	[0,4,6]	8000	[1,0.1,0.01]

Table 6: The settings of hyperparameters in multi-step scheduler in various datasets (see A.2).

Algorithm 1 Pseudo-code of object-level and frame-level prototype matching in a PyTorch-like style.

873 874 # q_obj_p: query object prototype (B_q x T x N x C) s_obj_p: support object prototype (B_s x T x N x C) 875 # q_frm_p: query frame prototype (B_q x T x C) 876 s_frm_p: support frame prototype (B_s x T x C) # 877 # B_q: number of query videos # B_s: number of support videos 878 # T: number of frames 879 # N: number of object prototypes in each frame 880 # OTAM: temporal metric 881 def matching(q_obj_p, s_obj_p, q_frm_p, s_frm_p):
 # Object-level prototype matching 882 883 q_obj_p = q_obj_p.reshape(B_q x T x N, C) # (B_q x T x N) x C s_obj_p = s_obj_p.reshape(B_s x T x N, C) # (B_s x T x N) x C 884 obj_sim = cos_sim(q_obj_p, s_obj_p) # (B_q x T x N) x (B_s x T x N) obj_dist = 1 - obj_sim # (B_q x T x N) x (B_s x T x N) 885 886 obj_dist = rearrange(obj_dist) # (B_q x T) x (B_s x T) x N x N 887 obj_fdist = obj_dist.min(3)[0].sum(2) + obj_dist.min(2)[0].sum(2) # (B_q x T) x (B_s x T) 888 obj_fdist = rearrange(obj_fdist) # B_q x B_s x T x T 889 obj_logits = obj_fdist.min(3)[0].mean(2) + obj_fdist.min(2)[0].mean 890 (2) # B_q x B_s 891 # Frame-level prototype matching 892 $q_frm_p = q_frm_p.reshape(B_q \times T, C) \# (B_q \times T) \times C$ s_frm_p = s_obj_p.reshape(B_s x T, C) # (B_s x T) x C 893 frm_sim = cos_sim(q_frm_p, s_frm_p) (B_q x T) x (B_s x T) 894 $frm_dist = 1 - frm_sim \# (B_q \times T) \times (B_s \times T)$ 895 frm_dist = rearrange(frm_dist) # B_q x B_s x T x T frm_logits = OTAM(frm_dist) # B_q x B_s 896 897 return obj_logits, frm_logits

C MORE EXPERIMENT RESULTS

C.1 IMPACT OF DIFFERENT NUMBERS OF PROMPTS

We conduct experiments on various prompt configurations in Table 7. First, we conduct experiments with 3, 6, and 12 spatial attributes. We can observe that the performance reaches its peak at G =6, possibly because too many spatial attributes may introduce local noise, while too few spatial attributes cannot afford enough spatial information. Furthermore, the best result is obtained when the number of temporal attributes is L = 3. We speculate that too few temporal attributes may fail to

Table 7: The ablation study on HMDB51 [31] to investigate the configuration of the number of spatial attributes G and temporal attributes L (see §C.1).

010				
913	(CI)	HMDB51		
914	$\{G,L\}$	1-shot	5-shot	
915	$\{3,3\}$	82.4	88.5	
916	$\{12, 3\}$	82.3	88.3	
917	$\{6, 1\}$	81.4	88.0	
	$\{6, 6\}$	81.5	88.0	
	$\{6,3\}$	82.6	88.7	

adequately convey the temporal changes in actions while too many temporal attributes often contain noisy temporal information, which leads to performance degradation. Therefore, we set G = 6 and L = 3 by default.

C.2 MORE VISUALIZATION EXAMPLES OF SPATIO-TEMPORAL ATTRIBUTES

We show more visualization results of spatial and temporal prompts in our DIST in Fig. 7. The attention maps of our DIST focus more on action-related objects and reduce attention to the back-ground and unrelated objects. This demonstrates our DIST grasps prior knowledge provided by spatial prompts to capture spatial details. Then, we calculate the cross-attention scores between temporal prompts and frames. It can be seen that different temporal prompts have different weights on the frame sequences, proving that our DIST can learn temporal relations and understand dynamic semantics. This further illustrates that our DIST can capture action-related spatial details and dynamic temporal information.



Figure 7: More visualization results of spatial and temporal prompts in our DIST under the 5-way 1shot setting. The spatial prompts are shown as highlighted response areas in each frame. Meanwhile, we show cross-attention temporal prompt weights in a line graph. See §C.2 for more details.

C.3 MODEL EFFICIENCY ANALYSIS

To analyze the effectiveness of training and inference, we list comparison results with SOTA CLIP FSAR [24] in terms of parameters, FLOPs, GPU memory, and inference speed in Table 8. We choose ViT-B as our visual encoder. As shown in Table 8, compared to SOTA CLIP-FSAR, our additional

Table 8: Complexity analysis for 5-way 1-shot HMDB51 [31] and Kinetics [33] evaluation. Here,
we report Params, FLOPs, GPU memory, and Speed for each model. "Acc¹" and "Acc²" are the
accuarcy on HMDB51 and Kinetics, respectively. See §C.3 for more details.

Method	Params	FLOPs	Memory	Speed	Acc^1	Acc^2
CLIP-FSAR	89.3M	901.9G	13.8G	36.7ms	75.8	89.7
DIST(ours)	97.2M	902.1G	14.1 G	40.9ms	82.6	92.7

Table 9: Some examples of action categories and their corresponding spatio-temporal decoupled knowledge generated by LLMs.

Category	Spatial	knowledge	Temporal knowledge
	1. swimming pool	2. diving board	1. Stand at the diving board.
dive	3. arms	4. legs	2. Jump into the water.
	5. pool	6. cliff	3. Resurface in the water.
	1. football	2. playground ball	1. Stand near a ball.
kick ball	3. feet	4. legs	2. Kick the ball.
	playground	6. park	3. The ball in motion, kick ended.
	1. bow	2. arrows	1. Bow, arrow, stance set.
archery	3. hands	4. arms	2. Drawing the bow, aiming.
	archery Range	outdoor Field	3. The arrow hits the target.
	1. toothbrush	2. toothpaste	 A person stands in front of the sink.
brush teeth	3. hands	4. mouth	2. Brushing teeth using a toothbrush and toothpaste.
	5. bathroom	6. bedroom	Rinsing mouth and cleaning toothbrush.
	1. eyeshadow palette	makeup brush	 Clean face, makeup tools ready.
apply eye makeup	3. hands	4. eye	2. Applying eye shadow, liner, and mascara.
	5. face	dressing Room	Finished eye makeup, enhanced eyes.
	 air drumming app 	2. drumsticks	 Tapping hands in rhythm.
air drumming	3. hands	4. arms	Mimicking drumming motions.
	5. bedroom	music studio	3. Slowing down, stopping drumming movements.
	1. candle	2. cake	1. Lit candles on a surface.
blowing out candles	3. hands	4. mouth	Blowing air to extinguish the flames.
	birthday party	cake shop	3. Candles are extinguished, and the celebration ends.
	 bowling ball 	bowling shoes	1. Holding the bowling ball, ready to throw.
bowling	3. hand	4. arm	2. Releasing the ball, rolling down the lane.
	5. bowling alley	6. entertainment center	3. Pins knocked down or remaining.

overhead is minimal. However, our DIST can bring 6.8% and 3.0% accuracy improvements on HMDB51 and Kinetics over CLIP-FSAR, respectively.

D ADDITIONAL EXAMPLES OF SPATIO-TEMPORAL KNOWLEDGE

In Table 9, we display some additional examples of action categories and their corresponding spatio temporal knowledge. This representative knowledge is automatically generated by LLMs.

¹⁰⁰⁶ E DISCUSSION

1008 E.1 LIMITATION ANALYSIS

1010 DIST relies on large language models to generate high-quality spatial and temporal prompts, which 1011 affect the final performance. In addition, relying solely on a fixed number of spatial and temporal 1012 prompts may not be optimal for every category. An adaptive approach, customizing the number of 1013 prompts for each category, would likely be more effective. In the future, we will also explore a more 1014 unified and effective way to inject spatial and temporal prompts into visual features.

1016 E.2 SOCIAL IMPACT

This work proposes a novel framework to leverage spatiotemporal-decoupled prior knowledge from LLM to compensate for visual features, so as to achieve more accurate few-shot matching. Our framework has demonstrated its effectiveness over multiple common benchmarks. On the positive side, the approach advances few-shot action recognition accuracy, and is valuable for real-world applications in the automated understanding of human actions from videos [73, 74], e.g., human-robotic interaction in elderly care facilities [74], behavior analysis for nursing procedures [73]. For potential negative social impact, our DIST struggles to understand human actions across varying domains, which is a common limitation shared by all few-shot action recognition methods [13, 16, 15]. Hence in the future, we will broaden the few-shot action recognition capability to generalize across varying domains.