Text as Any-Modality for Zero-Shot Classification by Consistent Prompt Tuning

Xiangyu Wu Nanjing University of Science and Technology Nanjing, China wxy_yyjhl@njust.edu.cn

Yang Yang*
Nanjing University of
Science and Technology
Nanjing, China
yyang@njust.edu.cn

Feng Yu Nanjing University of Science and Technology Nanjing, China hubaak@njust.edu.cn

Jianfeng Lu*
Nanjing University of
Science and Technology
Nanjing, China
lujf@njust.edu.cn

ABSTRACT

The integration of prompt tuning with multimodal learning has shown significant generalization abilities for various downstream tasks. Despite advancements, existing methods heavily depend on massive modality-specific labeled data (e.g., video, audio, and image), or are customized for a single modality. In this study, we present Text as Any-Modality by Consistent Prompt Tuning (TaAM-CPT), a scalable approach for constructing a general representation model toward unlimited modalities using solely text data. TaAM-CPT comprises modality prompt pools, text construction, and modality-aligned text encoders from pre-trained models, which allows for extending new modalities by simply adding prompt pools and modality-aligned text encoders. To harmonize the learning across different modalities, TaAM-CPT designs intra- and intermodal learning objectives, which can capture category details within modalities while maintaining semantic consistency across different modalities. Benefiting from its scalable architecture and pre-trained models, TaAM-CPT can be seamlessly extended to accommodate unlimited modalities. Remarkably, without any modality-specific labeled data, TaAM-CPT achieves leading results on diverse datasets spanning various modalities, including video classification, image classification, and audio classification. The code is available at https://github.com/Jinx630/TaAM-CPT.

CCS CONCEPTS

ullet Computing methodologies o Image representations.

KEYWORDS

Prompt Tuning, Multimodal Learning, Zero-shot Classification

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

MM '25, October 27-31, 2025, Dublin, Ireland

© 2025 Copyright held by the owner/author(s). Publication rights licensed to ACM. ACM ISBN 979-8-4007-2035-2/2025/10...\$15.00 https://doi.org/10.1145/3746027.3755288

ACM Reference Format:

Xiangyu Wu, Feng Yu, Yang Yang, and Jianfeng Lu. 2025. Text as Any-Modality for Zero-Shot Classification by Consistent Prompt Tuning. In *Proceedings of the 33rd ACM International Conference on Multimedia (MM '25), October 27–31, 2025, Dublin, Ireland.* ACM, New York, NY, USA, 15 pages. https://doi.org/10.1145/3746027.3755288

1 INTRODUCTION

As unified architectures [9, 15, 50] and multimodal pre-training models [8, 19, 37, 51] progress, recent works have exhibited impressive representation abilities in multimodal learning [13, 22, 59, 68, 72, 75]. In scenarios restricted by either labeled data or computational resources, owing to the aligned pre-trained models [37, 55, 61], prompt tuning [57, 60, 69] showcases robust generalization capabilities across various downstream tasks by adjusting a negligible number of parameters.

Current prompt tuning techniques rely heavily on massive modality-specific labeled data (e.g., video, audio, and image). For instance, as depicted in Figure 1 (a)(e), image supervised methods [18, 74] design text prompt to align with labeled image data for image-level tasks. Likewise, for video and audio level tasks, previous methods [20, 21, 29] adapt pre-trained models to downstream tasks supervised with labeled data. However, sufficient modality-specific labeled data necessitates considerable manual effort, which, in the face of labeled data limits, can impede the development of robust object classification networks. In the absence of labeled data altogether, these techniques may even fail outright.

To avoid above issue, some works advocate using the well-aligned embedding space, achieved by contrastive learning (e.g.,CLIP [6]), for prompt tuning. TAI-DPT [16], as a pioneering work depicted in Figure 1 (b)(e), proposes to enable labeled text data (e.g., cococaption [27]) instead of labeled image data for learning text prompt, while testing with images and learned text prompt. Similarly, PT-Text [26] pioneers the approach of audio-free prompt tuning, where the text prompt is learned from text rather than audio. To further reduce the manual cost of labeled text data, PVP [58] and TAI-Adapter [76] recommend using synthetic text data generated by LLMs [47] as a substitute. However, these strategies require the design of sophisticated text prompts, visual prompts, or adapter frameworks, as well as the deployment of a text encoder to encode

 $^{^{\}star}$ Corresponding author.

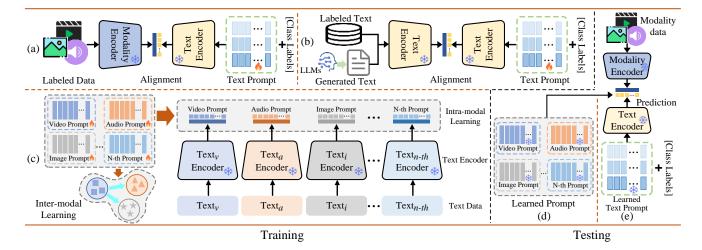


Figure 1: Different prompt tuning by frozen pre-trained encoders. (a)(b). Supervised methods with labeled and text data. (c). TaAM-CPT. Prompt tuning toward unlimited modalities without prompt encoding processes. (d). Testing of TaAM-CPT. (e). Testing of previous works.

the prompts. Additionally, these approaches focus solely on a single modality (e.g., video, image, or audio classification), and for more modalities, multiple independent models need to be trained additionally.

In this paper, we explore a universal representation model capable of scaling to unlimited modalities without any modality-specific labeled data. This necessitates the following conditions: 1) The model exclusively relies on easy-collected text data for training, eliminating the need for any labeled data. 2) The model architecture needs to be flexible enough to accommodate new categories or modalities and simplify the design of the prompt, thereby reducing the complexity of prompt encoding. 3) The model must ensure that learning across different modalities does not mutually affect each other, and appropriate training objectives should be designed to enhance the representational capabilities of all modalities.

Motivated by these factors, as shown in Figure 1 (c) and Figure 1 (d), we propose Text as Any-Modality for Consistent Prompt Tuning (TaAM-CPT), a general representation model toward unlimited modalities solely using text data generated by LLMs. Unlike TAI-DPT [16] and PVP [58], which require intricate, multi-grained text prompt designs, our method simplifies the design by characterizing any modality category as a randomly initialized vector. Leveraging the instruction following ability of LLMs [47], we can comfortably obtain text training data for any category. By directly optimizing the vectors within the aligned space of pre-trained models [37, 55, 61], we eliminate intermediate encoding processes. Since the initialization way for each category is identical, TaAM-CPT ensures the flexible addition of any category from any modality without retraining the already learned class-specific prompt. Moreover, we design a uni-directional contrastive loss, which uses modalities with stronger representational abilities to guide the learning of those weaker ones. Surprisingly, not only does it enhance the representational abilities of weaker modalities, but it also further improves the representational abilities of stronger modalities.

We conduct extensive experiments across multiple modalities and datasets, including video, audio, and image classification tasks. Without any labeled data, TaAM-CPT achieves superior performance to pre-trained models and recent SOTAs.

2 RELATED WORK

2.1 Video, Image, and Audio Classification.

Video classification involves identifying actions in the video. Early works [12, 48, 56] focus on designing two-stream networks and 3D CNNs for action recognition. Building on the success of transformers in the image, recent works [25, 53, 63, 64, 70, 71] explore effective objectives for adapting pre-trained image models to video understanding. To handle the problem of local video redundancy, UniFormerV2 [23] introduces local and global relation aggregators to learn discriminative representations.

Image classification aims to recognize all the categories in an image. To explore the correlations among labels, some works propose to incorporate semantic dependencies via object proposals [30, 54], semantic graph [73, 77], and transformer-based architecture [1, 41]. When labeled data is limited, another line of work [31, 32, 44] attempts to solve more challenging scenarios, including zero-shot, few-shot, and partial-label tasks.

Audio classification involves tagging audio signals into different categories. Traditional works [17, 34] mainly rely on machine learning technology and manual feature extraction. In recent years, driven by advancements in deep learning, some works [40, 62] have begun to explore the application of neural networks. Additionally, some efforts [14, 28] attempt to apply the transformer to audio classification, thereby capturing the long-term dependencies.

2.2 Prompt Tuning in Multimodal Learning.

Prompt tuning [10, 21, 52, 74] has emerged for rapidly adapting to downstream tasks by adjusting a minimal number of parameters.

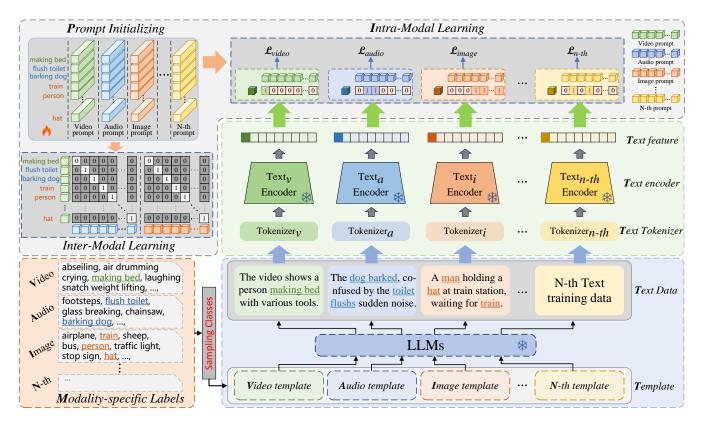


Figure 2: TaAM-CPT overview. We represent any category as a class-specific prompt and use LLMs to generate text data. Intra-modal learning aims to learn each prompt pool by pre-trained models. Inter-modal learning utilizes stronger modalities to guide those weaker ones.

For instance, some works [35, 74] introduce learnable context vectors to align with images via frozen CLIP encoders. When labeled data is limited, TAI-DPT [16] and PT-Text [26] introduced multigrained text prompts, surpassing pre-trained multimodal models in image and audio classification tasks, solely training text data. TAI-Adapter [76] combines LLM-driven data generation and cross-modal learning to enhance multi-label image classification tasks. PVP [58] further enhances image classification performance by co-learning pseudo-visual prompts and text prompts.

3 METHODS

The overview architecture of our proposed TaAM-CPT is illustrated in Figure 2. As shown, TaAM-CPT is designed as a general representation model toward unlimited modalities using only text data for prompt learning, which mainly consists of three parts: a) LLMs-assisted data construction, b) Prompt initializing and modality text encoding, and c) Intra- and inter-modal learning.

3.1 LLMs-Assisted Data Construction

Unlike noun filters used in TAI-DPT [16] and PVP [58], we construct prompt templates to instruct LLMs to generate text sentences that contain the given labels, as shown in Figure 2. For any given labels, we design the following query template:

TEMPLATE: Making several English sentences to describe a { **Modality** }. Requirements: Generate 5 English sentences! Each sentence should be less than 25 words and includes: { **Labels** }.

{ Modality } is populated with "video", "audio", "image", etc, and { Labels } indicates modality-specific labels, with a maximum of 2 for video modality, 3 for image and audio modalities. This design avoids the diversity (e.g., singular and plural) caused by noun filtering, and avoids noun filtering to process the phrases describing video and audio. Therefore, by generating text sentences containing these labels through LLMs, the corresponding ground truth for each sentence is from the { Labels } in the template. More details of prompt templates and text data generated by LLMs are provided in the appendix.

3.2 Prompt Initializing and Text Encoding

Prompt Initializing. We take video (\mathcal{V}), audio (\mathcal{A}), and image (\mathcal{I}) modalities as examples to introduce TaAM-CPT and demonstrate its potential for extension toward unlimited modalities. For each modality, we maintain a modality-specific prompt pool, defined as follows:

$$\mathbf{P}_{m} = [\mathbf{p}_{1}^{m}, \mathbf{p}_{2}^{m}, \mathbf{p}_{3}^{m}, ..., \mathbf{p}_{N}^{m}], \tag{1}$$

where $m \in \{\mathcal{V}, \mathcal{A}, I\}$ represents different modalities; $\boldsymbol{p}_i^m \in \mathbb{R}^d$ denotes *i*-th class-specific prompt; N denotes the total number

of labels. Note that the length of the prompt pool is identical for each modality (i.e., $\mathbf{P}_m \in \mathbb{R}^{N \times d}, m \in \{\mathcal{V}, \mathcal{A}, I\}$), encompassing all labels across all modalities. When a new modality emerges, a new modality-specific prompt pool will be created, avoiding affecting the already learned other prompt pools. When a new label arises, a new class-specific prompt will also be added to each prompt pool, avoiding affecting the existing class-specific prompts. Therefore, TaAM-CPT can be easily extended to unlimited modalities and categories.

Text Encoding. According to previous methods [16, 26, 58, 66], text is treated as a surrogate for other modalities for zero-shot classification. Such a paradigm potentially assumes that pre-trained models have aligned text with other modalities into a shared embedding space, thereby making it feasible to extract text features as substitutes for other modalities. However, these methods are designed for individual modalities and fail to utilize complementary information among multiple modalities. Hence, as shown in Figure 2, we adopt a parallel architecture and obtain modality-aligned text encoders (Text_v, Text_a, and Text_i Encoder) from pre-trained models ViCLIP, CLAP, and CLIP, to extract text features. Furthermore, we find CLIP and CLAP have superior representation abilities for image and audio, compared to ViCLIP for video, specifically reflected in the training loss and convergence speed of video modality. Inspired by the discovery, we design an uni-directional learning strategy to use stronger modalities to guide the learning of weaker modalities. We find that uni-directional learning can improve the performance for all modalities simultaneously.

3.3 Intra- and Inter-modal Learning

To learn the modality prompt pool for each modality, our work is to design two learning objectives: a) intra-modal learning aims to optimize the prompt pool for each modality using global text features. b) inter-modal learning aims to improve the representational abilities of weaker modalities based on stronger ones.

Intra-modal Learning. To make it easier, we take image modality as an example to introduce intra-modal learning, and the same approach is applied to video and audio modalities. The candidate label set is represented as $C = \{l_1, l_2, ..., l_N\}$, where N is the total number of labels across all modalities. Then, we denote the text training data for image labels as $\mathcal{T} = \{t_i, y_i\}_{i=1}^M$, where M is the number of texts; $y_i = \{y_{i1}, y_{i2}, ..., y_{i,N}\}$ denotes the ground truth of the text t_i and y_{ij} for $j \in \{1, 2, ..., N\}$ is 1 if the t_i is generated from the label l_j and 0 otherwise. Then, the text embedding of t_i is extracted by frozen text encoder of CLIP [6], formulated as follows:

$$\mathbf{h}_i = \phi(t_i),\tag{2}$$

where ϕ denotes the text encoder of CLIP, $\mathbf{h}_i \in \mathbb{R}^d$ denotes the normalized global text feature of t_i with d being the dimension. When processing the input text data of video or audio modalities, we simply replace ϕ as the text encoder of ViCLIP or CLAP to extract the corresponding text feature. The similarity of t_i and the prompt pool of image modality can then be computed by:

$$s_{ij} = \langle \mathbf{h}_i, \ \mathbf{p}_j \rangle, \quad \forall j \in \{1, 2, 3, ..., N\},$$
 (3)

where \mathbf{p}_j denotes the *j*-th prompt in the prompt pool of image modality. Note that the prompt can be optimized directly without

processing through any encoder or MLP. Such a paradigm simplifies the design of the prompt and reduces the computational cost by half. For the optimization of the prompt, we employ Ranking loss instead of InfoNCE or Cross-Entropy loss, since InfoNCE loss requires massive negative samples and high-cost softmax function to optimize well, Cross-Entropy loss only optimizes positive labels while ignoring loss from negative labels, leading to very slow convergence. Ranking loss is computed by:

$$\mathcal{L}_{\mathbf{I}} = \frac{1}{B} \sum_{k=1}^{B} \sum_{i \in \{c^{+}\}} \sum_{j \in \{c^{-}\}} \max(0, m - s_{ki} + s_{kj}), \tag{4}$$

where c^+ denotes positive labels with y_{ij} for $j \in \{1, 2, ..., N\}$ is 1, c^- denotes negative labels, s_{ki} and s_{kj} are positive pair and negative pair similarities described in Eq. (3), m is denoted as the margin to measure the difference between each pair of similarities. For the video and audio modalities, we substitute the text encoder ϕ described in Eq. (3) to the text encoder of ViCLIP and CLAP to obtain the text feature, and then compute the similarities between the text feature and video prompt pool, audio prompt pool. As a result, we can obtain the Ranking loss \mathcal{L}_V and \mathcal{L}_A and \mathcal{L}_I . The total loss for intra-modal learning can be written as:

$$\mathcal{L}_{intra} = \mathcal{L}_{I} + \mathcal{L}_{V} + \mathcal{L}_{A}. \tag{5}$$

During training, we fix text encoders and optimize the modality-specific prompt pools by Eq. (5). Note that the positive labels in Eq. (4) only contain positive image labels, while negative labels contain not only negative image labels but also labels from other modalities. Other modalities' labels serving as negative labels not only expand the number of negative pairs but also enhance the representational ability of the video modality. By analogy, this rule can be applied to audio and image modalities also.

Inter-modal Learning. The discrepancy in the information content of image, audio, and video modalities results in a significant modality gap between the aligned video and text modalities and subpar zero-shot classification performance. Motivated by this phenomenon, we propose uni-directional contrastive learning, which guides the learning of weaker modalities using the stronger ones. In this paper, we adaptively determine the weak modality during training based on the lowest validation performance. Specifically, the video modality is treated as weak as its performance is always lower, and image and audio as stronger ones. To facilitate understanding, we rephrase Eq. (1) into the follow format:

$$\begin{split} \mathbf{P}_{\mathcal{V}} &= [\ \mathbf{p}_{v_{1}}^{\mathcal{V}}, \mathbf{p}_{v_{2}}^{\mathcal{V}}, ..., \mathbf{p}_{v_{v}}^{\mathcal{V}}, \mathbf{p}_{a_{1}}^{\mathcal{V}}, \mathbf{p}_{a_{2}}^{\mathcal{V}}, ..., \mathbf{p}_{a_{d}}^{\mathcal{V}}, \mathbf{p}_{w_{1}}^{\mathcal{V}}, \mathbf{p}_{w_{2}}^{\mathcal{V}}, ..., \mathbf{p}_{w_{w}}^{\mathcal{V}} \], \\ \mathbf{P}_{\mathcal{A}} &= [\ \mathbf{p}_{v_{1}}^{\mathcal{A}}, \mathbf{p}_{v_{2}}^{\mathcal{A}}, ..., \mathbf{p}_{v_{v}}^{\mathcal{A}}, \mathbf{p}_{a_{1}}^{\mathcal{A}}, \mathbf{p}_{a_{2}}^{\mathcal{A}}, ..., \mathbf{p}_{a_{d}}^{\mathcal{A}}, \mathbf{p}_{w_{1}}^{\mathcal{A}}, \mathbf{p}_{w_{2}}^{\mathcal{A}}, ..., \mathbf{p}_{w_{w}}^{\mathcal{A}} \], \\ \mathbf{P}_{I} &= [\ \mathbf{p}_{v_{1}}^{\mathcal{I}}, \mathbf{p}_{v_{2}}^{\mathcal{I}}, ..., \mathbf{p}_{v_{v}}^{\mathcal{I}}, \mathbf{p}_{a_{1}}^{\mathcal{I}}, \mathbf{p}_{a_{2}}^{\mathcal{I}}, ..., \mathbf{p}_{a_{d}}^{\mathcal{A}}, \mathbf{p}_{w_{1}}^{\mathcal{I}}, \mathbf{p}_{w_{2}}^{\mathcal{I}}, ..., \mathbf{p}_{w_{w}}^{\mathcal{I}} \], \end{split} \tag{6}$$

where v+a+w=N, $\mathbf{p}_k^{\mathcal{V}}$, $\mathbf{p}_k^{\mathcal{H}}$ and $\mathbf{p}_k^{\mathcal{I}}$ represent class-specific prompt of video, audio, and image prompt pools. Note that the initialized prompt pool of each modality is identical, which means the prompt pool of the video modality contains video labels of size v, audio labels of size a, and image labels of size w. The prompt pool for image and audio modalities is the same as the video modality.

We then present the uni-directional contrastive objective based on $\mathbf{P}_{\mathcal{V}}$ and $\mathbf{P}_{\mathcal{A}}$. Specifically, the similarity matrix can be computed by $\mathbf{P}_{\mathcal{A}}^{\mathsf{T}}\mathbf{P}_{\mathcal{V}} \in \mathbb{R}^{N \times N}$. And the ground truth for $\mathbf{P}_{\mathcal{V}}$ and $\mathbf{P}_{\mathcal{A}}$ of N labels is a diagonal matrix. Note that the size of the similarity matrix and ground truth matrix is batch-size agnostic and equals

the number of total labels. Therefore, for each video prompt of $P_{\mathcal{V}}$ and audio prompt of $P_{\mathcal{A}}$, the softmax-normalized video prompt to audio prompt and ground truth matrix can be defined as:

$$p_{ij}^{v2a} = \frac{\exp(s(v_i, a_j)/\tau)}{\sum_{k=1}^{N} \exp(s(v_i, a_k)/\tau)}, \quad y^{v2a} = \begin{bmatrix} 0 & 0 & \cdots & 0 & \cdots & 0 \\ 0 & \ddots & \ddots & \vdots & & \vdots \\ \vdots & \ddots & 0 & 0 & \cdots & 0 \\ 0 & \cdots & 0 & 1 & \cdots & 0 \\ \vdots & & \vdots & \vdots & \ddots & \vdots \\ 0 & \cdots & 0 & 0 & \cdots & 1 \end{bmatrix}, (7)$$

where v_i and a_j denote video prompt and audio prompt, $s(\cdot,\cdot)$ represents similarity function. Note that the ground truth label y^{v2a} is different from the label matrix in vanilla contrastive learning (i.e. identity matrix), where the first v+w diagonal elements are set to 0. It indicates that the loss generated at these positions will be ignored when calculating the cross-entropy loss. Therefore, the unidirectional contrastive loss for P_V and $P_{\mathcal{A}}$ can be defined as $\mathcal{L}_{v2a} = \mathcal{L}_{CE}(y^{v2a}, p^{v2a})$, where $y_{ij}^{v2a} \in \{0,1\}$ for $\forall i,j \in \{1,2,\ldots,N\}$ represents the similarity ground truth between video prompt v_i and audio prompt a_j . Similarly, we can obtain the uni-directional contrastive loss between the prompt pool of video modality P_V and prompt pool of image modality P_I : $\mathcal{L}_{v2w} = \mathcal{L}_{CE}(y^{v2w}, p^{v2w})$. And the total inter-learning loss can be defined as:

$$\mathcal{L}_{inter} = \mathcal{L}_{v2a} + \mathcal{L}_{v2w}. \tag{8}$$

Consequently, we align the prompts of image labels in the video prompt pool with those in the image prompt pool, and the prompts of audio labels with those in the audio prompt pool. These aligned image and audio prompts will be treated as negative labels for training video prompt pool in intra-modal learning, thereby expanding the number of negative pairs. In addition, diagonal elements corresponding to video and image labels in the ground truth matrix are set to 0, which avoids affecting the learning of the prompt of video labels. During training, we apply uni-directional contrastive learning to video-to-audio and video-to-image. The total loss of TaAM-CPT is: $\mathcal{L}_{total} = \lambda_1 \mathcal{L}_{intra} + \lambda_2 \mathcal{L}_{inter}$, where λ_1 and λ_2 denote the loss weights of intra-modal learning and inter-modal learning.

3.4 Discussion

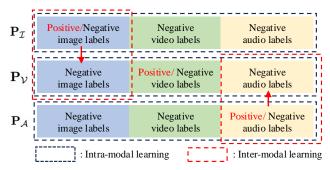


Figure 3: Visualization of learning process.

In this subsection, as illustrated in Figure 3, we discuss how intra- and inter-modal learning work well. For inter-modal learning, we employ uni-directional contrastive learning, aligning "negative image/audio labels" from the "video prompt pool" with "positive/negative image labels" from the "image prompt pool" and

"positive/negative audio labels" from the "audio prompt pool". We can actually treat this process as transferring knowledge from the "image/audio prompt pool" to the "video prompt pool". For intramodal learning, take the image modality as an example. Although the "negative labels" contain "negative image/video/audio labels" in the "image prompt pool", these "negative video/audio labels" don't require modality alignment. The purpose is just to increase the number of negative samples, thereby learning more robust representations of "positive image labels". For the video modality, the "negative labels" come from aligned "negative image/audio labels" and "negative video labels" in the "video prompt pool". Such a uni-directional contrastive learning strategy ensures that "negative image/audio labels" in the "video prompt pool" can not affect the learning of "positive image/audio labels" in the "image/audio prompt pool".

3.5 Model Testing

After learning the prompt pool of each modality, each prompt uniquely represents a specific class. We take the video modality as an example to showcase the video classification. Given an input video, we replace the video modality-specific text encoder in Figure 2 with the video encoder of ViCLIP to obtain the video feature. Then, we directly calculate the similarity between the video feature and each prompt in the video prompt pool, and the prediction of the input video is the prompt with the highest similarity. It can be seen that each prompt is calculated directly with the video features without any encoding processing, which significantly improves the inference speed of the model. For image classification and audio classification, we adopt the same approach, calculating the similarity between image or audio and their corresponding prompt pools to obtain the predictions.

4 EXPERIMENTS

4.1 Experimental Setup

Datasets. We conduct extensive experiments on 13 datasets across video, image, and audio modalities. For video classification, we select UCF101 and the large-scale datasets Kinetic-400/600/700. For image classification, besides MSCOCO, VOC2007, and NUSWIDE used in previous works, we also select the VOC2012, ImageNet-mini, and Objects365 to evaluate our method. For audio classification, we follow PT-Text, selecting ESC50 and US8K.

Implementation Details. We select the pre-trained models, open-sourced by the LAION [42], as the modality-specific encoders, i.e., ViCLIP-Base for video modality, CLIP-ViT-B-32 for image modality, and CLAP for audio modality. The LLaMA-2-7B is selected for generating 100k text sentences for each modality, on a single Tesla V100, it takes about 2 hours. By simply adding some spatial relationships instructions in the template, LLaMA-2-7B can generate texts that accurately reflect spatial relationships. For each class-specific prompt, we initialize it as a vector with a length of 512, mean being 0, and std being 0.02. During training, all modality-aligned text encoders are fixed, and only prompts are optimized. We evaluate TaAM-CPT by top-1/5 accuracy and mean average precision (mAP) metrics. See appendix for more implementation details.

Table 1: Comparison with ZS-CLIP and SOTAs on zero-shot image classification.

Methods MSCOCO VOC2007 VOC2012 NUSWIDE ImageNet-mini Objects365 ZS - CLIP_[ICLR24] (85.5, 94.3) 55.6 80.5 80.1 37.1 19.8 TAI - DPT_[CVPR23] 65.1 88.3 85.1 (86.2, 94.7) 24.1 46.5 TAI - Adapter [arXiv23] 85.5 53.3 (86.7,94.4) 25.8 67.7 89.0 Data - free [arXiv24] 66.8 88 7 86.0 47.0 (86.1, 94.9) 23 9 $PVP_{\hbox{\tt [IJCAI24]}}$ 49.3 (87.4, 95.3)67.7 88.9 86.2 26.3 TaAM - CPT(Ours) 87.8 49.6 28.2 68.1 89.4 (90.4, 98.3)

Table 2: Comparison with ZS-CLAP and SOTAs on zeroshot audio classification.

Methods	ESC50	US8K
ZS - CLAP _[ICASSP23]	90.5	<u>76.2</u>
PT - Text _[ICASSP24]	93.9	_
TaAM - CPT(Ours)	94.2	85.2

Table 3: Results with ZS-ViCLIP on zero-shot video classification.

	UCI				K6		K7	
Wiethous	top1	top5	top1	top5	top1	top5	top1	top5
ZS-ViCLIP _[ICLR24] TaAM-CPT(Ours)	73.3 75.4	93.3 95.7	53.8 55.2	78.7 80.4	52.0 52.9	78.4 80.1	43.5 46.0	68.6 71.1

4.2 Results on Zero-Shot Tasks

To evaluate TaAM-CPT, besides the zero-shot performance comparison with pre-trained multimodal models (*i.e.* ViCLIP, CLIP, CLAP), we also compare its performance with existing SOTA methods on image classification and audio classification tasks. Notably, in the zero-shot video classification field, there has been no research that explores a similar training setting, *i.e.*, solely training with text data for prompt tuning. Therefore, we only select ViCLIP as the zero-shot benchmark for comparison.

Video Classification. We adopt the default prompt "a video of a [CLASS]" to obtain the zero-shot results of ViCLIP. From Table 3, TaAM-CPT outperforms ZS-ViCLIP by 2.1% top-1 and 2.4% top-5 accuracy on UCF101. On the larger Kinetic-400/600/700 datasets, respectively, TaAM-CPT also surpasses ZS-ViCLIP by 0.9~3.0% top-1 and top-5 accuracy on all datasets, showing the effectiveness of TaAM-CPT without labeled video data.

Image Classification. For zero-shot image classification, we present the results in Table 1 and compare our approach with SO-TAs TAI-DPT [16], TAI-Adapter [76], Data-free [66], and PVP [58] trained with complex prompt design or adapter module. The results of ZS-CLIP are obtained by inputting the default prompt "a photo of a [CLASS]" to CLIP. From Table 1, TaAM-CPT outperforms ZS-CLIP by a large margin of 12.5% and 12.4% mAP on MSCOCO and NUSWIDE, respectively. On VOC2007 and VOC2012, TaAM-CPT also improves by 7.0% ~9.0% over ZS-CLIP.

Audio Classification. The results for zero-shot audio classification with CLAP and recent SOTA PT-Text [26] are shown in Table 2. Our TaAM-CPT outperforms ZS-CLAP with 3.7% and 9.0% accuracy on ESC50 and US8K, despite the high performance of CLAP. Furthermore, without intricate prompt design, TaAM-CPT surpasses PT-Text 0.3% on the ESC50 dataset.

4.3 Integrating with other methods

Following TAI-DPT [16], we conduct the experiments of integrating TaAM-CPT with other supervised models in an off-the-shelf manner, further improving their performance. Take a video with n labels as an example, the supervised model's softmax predictions denote as $P_S = (p_{s1}, p_{s2}, ..., p_{sn})$. For TaAM-CPT, we calculate the similarity between video and n class-specific video prompts and obtain softmax predictions $P_T = (p_{t1}, p_{t2}, ..., p_{tn})$. Therefore, the integrated results can be computed by $P_I = (p_{s1} + p_{t1}, p_{s2} + p_{t2}, ..., p_{sn} + p_{tn})$.

Table 4: Integrating TaAM-CPT with supervised video models.

Methods	K4	100	Ke	K600		K700	
Methods	top1	top5	top1	top5	top1	top5	
Video Swin _[CVPR22]	82.7	95.5	84.0	96.5	-	_	
+TaAM-CPT(Ours)	83.5	95.9	84.8	97.1	_	_	
MTV _[CVPR22]	81.8	95.0	83.8	96.2	73.5	90.3	
+TaAM-CPT(Ours)	82.9	95.7	84.7	97.0	74.8	91.2	
AIM _[ICLR23]	83.9	96.3	-	-	76.9	92.1	
+TaAM-CPT(Ours)	84.6	97.2	-	_	77.2	93.0	
UniFormerV2 _[ICCV23]	84.0	96.3	84.8	96.8	75.4	92.6	
+TaAM-CPT(Ours)	84.8	97.1	85.5	97.6	76.1	93.4	
UMT _[ICCV23]	85.7	97.0	87.8	97.8	78.5	94.3	
+TaAM-CPT(Ours)	86.2	97.6	88.1	98.0	78.8	94.7	

Video Classification. We select the Base size model of Video Swin Transformer [33], MTV [65], AIM [67], UniFormerV2 [23], and UMT [24] as baselines. The results are shown in Table 4. After integrating our TaAM-CPT with Video Swin, MTV, AIM-B, UniFormerV2-B, and UMT-B on Kinetic-400/600/700 datasets, the video classification performance of these methods can be further improved, while these methods achieve promising performances.

Image Classification. In Table 6, we select the newest DualCoOp++ [18] instead of DualCoOp [46] used in previous SO-TAS [16,58], and reproduce DualCoOp++ on these datasets (marked with *). + indicates integrating predictions with DualCoOp++*. In Table 6, while DualCoOp++* obtains promising performance, +TaAM-CPT can further enhance the image classification results. Compared with +TAI-DPT and +PVP, our +TaAM-CPT achieves higher performance in all cases, and surpasses +PVP by considerable margins of 0.2%, 0.3%, and 1.2% mAP on these datasets. Notably, TAI-DPT and PVP rely on costly prompt encoders and are only customized for a single image modality. Our TaAM-CPT is a general representation model that can accommodate unlimited modalities and class labels.

Table 5: Integrating TaAM-CPT with supervised audio methods.

Methods	ESC50	US8K
HTS-AT _[ICASSP22]	97.0	94.7
+ TaAM-CPT (Ours)	97.2	95.1
CrissCross[AAAI23]	90.5	92.1
+TaAM-CPT(Ours)	94.7	92.8

Audio Classification. We also study the audio classification results of integrating with HTS-AT [5] and CrissCross [40]. In the same video classification task, we compute the similarities between the audio feature and the audio prompt pool as the predictions.

Methods 10% 20% 30% 40% 50% 60% 70% 80% 90% Avg $DualCoOp_{\hbox{\tt [NeurIPS22]}}$ 81.0 82.3 82.9 83 4 83 5 83 9 84 0 84 1 84 3 833 DualCoOp++[TPAMI24] 81.4 83.1 83.7 84.2 84.4 84.5 84.8 85.0 85.1 84.0 **MSCOCO** DualCoOp++*[TPAMI24] 81.5 83.2 84.084.4 84.5 84.7 84.8 85.1 85.2 84.1 +TAI-DPT_[CVPR23] 83.3 84.5 84.5 84.7 85.0 85.1 85.2 85.2 84.3 81.7 +PVP[I]CAI24] 82.1 83.6 84.5 84.7 85.0 85.3 85.3 85.6 85.6 84.6 +TaAM-CPT(Ours) 82.4 83.8 84.6 85.0 85.1 85.3 85.5 85.7 85.8 84.8 93.8 93.8 94.3 94.7 94.8 94.9 DualCoOp[NeurIPS22] 91.4 94.6 94.9 94.1 $DualCoOp++_{\hbox{\tt [TPAMI24]}}$ 92.7 93.4 93.8 94.0 94.3 94.7 94.9 94.4 94.4 94.1 DualCoOp++*[TPAMI24] 93.0 94.4 94.9 93.9 94.2 94.6 94.8 95.1 95.0 94.4 $+ {\rm TAI\text{-}DPT}_{{\color{blue}[{\rm CVPR23}]}}$ 93.2 94.0 94.2 94.6 94.7 94.8 95.0 95.1 95.1 94.5 +PVP[IJCAI24] 93.5 94.3 94.4 94.6 95.0 95.1 95.2 95.2 95.3 94.7 +TaAM-CPT(Ours) 93.9 94.6 94.8 95.1 95.3 95.4 95.4 95.5 95.6 95.0 56.2 56.9 57.4 57.9 57.9 57.6 58.2 58.8 DualCoOp[NeurIPS22] 54.0 57.2 NUSWIDE $DualCoOp + + *_{\hbox{$\tt [TPAMI24]}}$ 54.4 56.6 58.1 58.7 58.9 59.3 59.7 59.8 60.1 58.4 +TAI-DPT_[CVPR23] 56.9 58.1 58.5 58.8 58.8 59.1 59.1 59.5 60.0 58.7 +PVP[I]CAI24] 57.3 58.6 59.3 59.4 59.6 60.0 60.1 60.1 60.3 59.4

60.7

60.8

61.3

Table 6: The mAP results for partial-label setting on all datasets, where the performance of +TAI-DPT/+PVP/+TaAM-CPT integrates the predictions of TAI-DPT/PVP/TaAM-CPT and DualCoOp++*.

From Table 5, the performance of both HST-AT and CrissCross is enhanced on ESC50 [36] and US8K [39] datasets.

58.2

59.6

60.5

+TaAM-CPT(Ours)

Further Analysis

We conduct further analysis to explore TaAM-CPT. More results (e.g., each component, hyperparameter, prompt dimension, more datasets, training data size, etc.) are presented in the appendix.

Quantity of modalities and categories. We first explore the feasibility of TaAM-CPT for unlimited modalities and categories. For clarity, we adopt Bert-base [8] with 110MB parameters as reference. TaAM-CPT initializes each prompt with a 512-d vector, meaning one class prompt occupies 512 parameters. Therefore, for Nmodalities, TaAM-CPT can accommodate approximately $\frac{110,000,000}{512\pi}$ class prompts. For example, for 10 modalities, the prompt pool size can reach 21484, which is sufficient to cover the common categories.

Table 7: Results of evaluating the unified architecture.

VP	IΡ	AP	\mathcal{L}_{Ia}	$ \mathcal{L}_{\mathrm{Ie}} $	K400	MSCOCO	ESC50
ZS-	-ViC	LIP,C	LIP,C	LAP	(53.8, 78.7)	55.6	90.5
\checkmark	×	×	✓	×	(53.8, 78.9)	_	-
×	✓	×	✓	×	_	65.8	-
×	×	✓	✓	×	_	_	92.5
\checkmark	✓	✓	✓	×	(53.7, 79.1)	65.2	92.7
\checkmark	✓	✓	✓	✓	(55.2, 80.4)	68.1	94.2

Unified Architecture. Our TaAM-CPT is designed as a general model toward unlimited modalities, exhibiting more robust object recognition capabilities than single modality-specific models. Table 7 presents the results of training each modality independently by intra-modal learning (e.g. $VP \checkmark with \mathcal{L}_{Ia} \checkmark$), as well as the impact of applying the uni-directional contrastive learning (\mathcal{L}_{Ie}) across modalities. We can see that training a single modality prompt by intra-modal learning has already yielded better results than the pre-trained models, and when all modalities are trained together, the performance of each modality can be further improved.

Table 8: Results of different learning manners.

61.3

61.7

60.6

61.4

$\mathcal{L}_{ ext{Inter}}$	K400	MSCOCO	ESC50
ZS-ViCLIP,CLIP,CLAP	(53.8, 78.7)	55.6	90.5
	(52.1, 79.3)	64.8	91.7
	(51.9, 79.5)	65.1	91.8
	(53.6, 79.4)	67.1	92.4
	(53.2, 79.2)	65.3	93.2
	(54.3, 79.8)	67.1	92.9
	(55.2, 80.4)	68.1	94.2

Inter-moda Learning. In Table 8, $\langle a, b \rangle \longrightarrow \langle c \rangle$ denotes unidirectional contrastive learning from a, b to c, while ←→ denotes naive bi-directional learning. Both $\langle I, V \rangle \longrightarrow \langle A \rangle$ and $\langle A, V \rangle \longrightarrow \langle I \rangle$ improve the performance of image and audio modalities while decreasing on video modality. Notably, $\langle I \rangle \longrightarrow \langle V \rangle$ and $\langle A \rangle \longrightarrow \langle V \rangle$ significantly outperform ZS-CLIP and ZS-CLAP by a large margin, demonstrating the effectiveness of inter-modal learning. Additionally, uni-directional learning can achieve higher performance than bi-directional learning on all datasets.

Table 9: Time complexity of adding new concepts compared to initial training time. IT: Initial Training, CT: Continue Training, NT: Novel Training, Fro: Frozen.

K400_Normal_30	T(58.1, 82.3)	CT(58.4, 82.5)	Fro(58.1 , 82.3)
MSCOCO_Normal_30	IT (65.4)	CT (65.6)	Fro (65.4)
K400_Novel_20	IT (56.3)	NT (56.7)	NT (56.5)
MSCOCO_Novel_20	IT (69.2)	NT (68.8)	NT (69.3)
ESC50_Novel_20	IT (92.1)	NT (92.5)	NT (93.0)
Training time	28min	30min	17min

Time Complexity. The trainable parameters of TaAM-CPT are the class-specific prompts. Therefore, for simplicity, we randomly sample K400 with 30 video labels and MSCOCO with 30 image

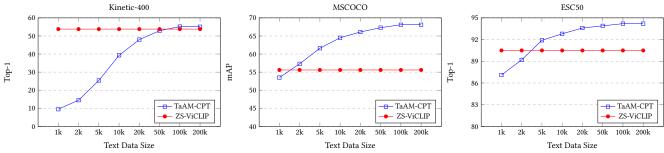


Figure 4: Results of different size of text training data on Kinetic-400, MSCOCO and ESC50 datasets.

labels as the initial label set. The novel concepts consist of additional labels (K400 with 20 video labels, MSCOCO with 20 image labels), and the novel audio modality (ESC50 with 20 labels), with each label containing 500 text descriptions using LLMs. The comparison of training time results is shown in Table 9, Among them, (K400_Novel_20, MSCOCO_Novel_20, and ESC50_Novel_20) represent the added new labels and new modality. We can see that for the normal labels (K400_Normal_30, MSCOCO_Normal_30), whether they continue training (Continue training) or are frozen (Frozen), neither affects the learning of the new modality labels and maintains the performance that matches the initial joint training (Initial training). For the novel modalities and labels, their Novel training performances are still comparable to Initial training.

Text Training Data Size. Our TaAM-CPT is trained with text data generated by LLMs instead of modality-specific labeled data. Therefore, we conduct various experiments with different sizes of text training data on the Kinetic-400, MSCOCO, and ESC50 datasets. As shown in Figure 4, on the Kinetic-400 dataset with text data size being 1k, the top-1 accuracy is only 9.8% due to the insufficient number of text data for each class, which hinders the learning of robust class-specific representations. However, as we continue to expand the scale of text training data, the corresponding text data for each class also increases gradually. When the text data reaches 100K, our TaAM-CPT outperforms ZS-ViCLIP. On the MSCOCO and ESC50 datasets, which contain 80 and 50 class labels, respectively, when the amount of text data is 5K, our method has already significantly surpassed ZS-CLIP and ZS-CLAP by 7% mAP and 2% top-1 accuracy. The performance on these two datasets begins to stabilize when the amount of text data is increased to 50K, indicating that datasets with more classes require a larger scale of text training data.

4.5 Visualization

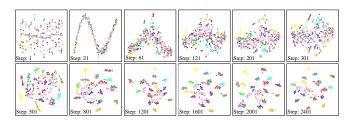


Figure 5: Distribution of video prompt and video feature.

Intra-modal Learning. We randomly selected 20 video classes on Kinetic-400. For each video sample, we computed its similarity

with each video prompt, resulting in a 400-d vector and using t-SNE [49] for visualization, which reflects the learning process of each video class prompt in Figure 5. Since the initialization method is identical, video samples from the same category show a uniform distribution before model training (*Step: 0*). As training progresses, the class-specific prompt begins to learn the unique representations (*Step: 21~1201*) for each category (*Step: 1601~2401*). Visualizations for more datasets can be found in the appendix.

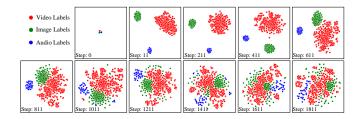


Figure 6: Distribution of prompt for different modalities.

Inter-modal Learning. We select Kinetic-400, MSCOCO, and ESC50 datasets, which contain 400, 80, and 50 class labels, respectively. As shown in Figure 6, before model training (*Step: 0*), the prompt pools for each modality are initialized in the same vector. When starting training, the distribution of different modalities rapidly separates (*Step: 11~211*), as each modality-first learns modality-specific representations through modality-aligned text encoders. As training progresses, uni-directional contrastive learning gradually pulls the representation space of the video modality towards image and audio modalities (*Step: 411~1411*), indicating that the video modality is continuously learning the representations of image and audio modalities. Furthermore, each modality still maintains its own representation space without being disrupted by the other modalities (*Step: 1611~1811*).

5 CONCLUSION

In this paper, we explore a scalable way of constructing a universal representation model for various modalities. Based on a flexible architecture and aligned pre-trained models, we develop TaAM-CPT, treating any category as a learnable vector and optimizing it directly through aligned pre-trained models. In addition, uni-directional contrastive learning also improves the classification performance of all modalities. The experimental results on 13 datasets show that TaAM-CPT achieves leading results in various classification tasks, including zero-shot video classification, image classification, audio classification, and partial-label image classification.

ACKNOWLEDGMENTS

This work is supported by the National Key RD Program of China (2022YFF0712100), NSFC (62276131), Natural Science Foundation of Jiangsu Province of China under Grant (BK20240081), the Fundamental Research Funds for the Central Universities (No.30922010317, No.30923011007)

REFERENCES

- Uzair Aslam Bhatti, Mengxing Huang, Harold Neira-Molina, Shah Marjan, Mehmood Baryalai, Hao Tang, Guilu Wu, and Sibghat Ullah Bazai. 2023. MFFCG– Multi feature fusion for hyperspectral image classification using graph attention network. ESWA 229 (2023), 120496.
- [2] Joao Carreira, Eric Noland, Andras Banki-Horvath, Chloe Hillier, and Andrew Zisserman. 2018. A short note about kinetics-600. CoRR (2018).
- [3] Joao Carreira, Eric Noland, Chloe Hillier, and Andrew Zisserman. 2019. A short note on the kinetics-700 human action dataset. CoRR (2019).
- [4] Joao Carreira and Andrew Zisserman. 2017. Quo vadis, action recognition? a new model and the kinetics dataset. In CVPR. 6299–6308.
- [5] Ke Chen, Xingjian Du, Bilei Zhu, Zejun Ma, Taylor Berg-Kirkpatrick, and Shlomo Dubnov. 2022. HTS-AT: A Hierarchical Token-Semantic Audio Transformer for Sound Classification and Detection. In ICASSP. 646-650.
- [6] Mehdi Cherti, Romain Beaumont, Ross Wightman, Mitchell Wortsman, Gabriel Ilharco, Cade Gordon, Christoph Schuhmann, Ludwig Schmidt, and Jenia Jitsev. 2023. Reproducible scaling laws for contrastive language-image learning. In CVPR. 2818–2829.
- [7] Tat-Seng Chua, Jinhui Tang, Richang Hong, Haojie Li, Zhiping Luo, and Yantao Zheng. 2009. NUS-WIDE: a real-world web image database from National University of Singapore. In CIVR.
- [8] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In NAACL. 4171–4186.
- [9] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. 2021. An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale. In ICLR.
- [10] Haoyi Duan, Yan Xia, Zhou Mingze, Li Tang, Jieming Zhu, and Zhou Zhao. 2024. Cross-modal Prompts: Adapting Large Pre-trained Models for Audio-Visual Downstream Tasks. NeurIPS 36 (2024).
- [11] Mark Everingham, Luc Van Gool, Christopher K. I. Williams, John M. Winn, and Andrew Zisserman. 2010. The Pascal Visual Object Classes VOC Challenge. IJCV 88, 2 (2010), 303–338.
- [12] Christoph Feichtenhofer. 2020. X3D: Expanding Architectures for Efficient Video Recognition. In CVPR. 200–210.
- [13] Zhongtian Fu, Kefei Song, Luping Zhou, and Yang Yang. 2024. Noise-Aware Image Captioning with Progressively Exploring Mismatched Words. In AAAI. 12091–12099.
- [14] Shefali Garg, Zhouyuan Huo, Khe Chai Sim, Suzan Schwartz, Mason Chua, Alëna Aksënova, Tsendsuren Munkhdalai, Levi King, Darryl Wright, Zion Mengesha, et al. 2024. Improving Speech Recognition for African American English with Audio Classification. In ICASSP. 12356–12360.
- [15] Yuan Gong, Yu-An Chung, and James R. Glass. 2021. AST: Audio Spectrogram Transformer. In ISCA. 571–575.
- [16] Zixian Guo, Bowen Dong, Zhilong Ji, Jinfeng Bai, Yiwen Guo, and Wangmeng Zuo. 2023. Texts as Images in Prompt Tuning for Multi-Label Image Recognition. In CVPR. 2808–2817.
- [17] Mikael Henaff, Kevin Jarrett, Koray Kavukcuoglu, and Yann LeCun. 2011. Unsupervised learning of sparse features for scalable audio classification. ISMIR (Jan 2011).
- [18] Ping Hu, Ximeng Sun, Stan Sclaroff, and Kate Saenko. 2023. Dualcoop++: Fast and effective adaptation to multi-label recognition with limited annotations. TPAMI (2023).
- [19] Longfei Huang, Xiangyu Wu, Jingyuan Wang, Weili Guo, and Yang Yang. 2024. Refining Visual Perception for Decoration Display: A Self-Enhanced Deep Captioning Model. In ACML, Vol. 260. 527–542.
- [20] Saksham Singh Kushwaha and Magdalena Fuentes. 2023. A multimodal prototypical approach for unsupervised sound classification. In INTERSPEECH. 266–270.
- [21] Bing Li, Jiaxin Chen, Xiuguo Bao, and Di Huang. 2023. Compressed Video Prompt Tuning. NeurIPS 36 (2023), 31895–31907.
- [22] Junnan Li, Dongxu Li, Silvio Savarese, and Steven C. H. Hoi. 2023. BLIP-2: Bootstrapping Language-Image Pre-training with Frozen Image Encoders and Large Language Models. In ICML, Vol. 202. 19730–19742.
- [23] Kunchang Li, Yali Wang, Yinan He, Yizhuo Li, Yi Wang, Limin Wang, and Yu Qiao. 2022. UniFormerV2: Spatiotemporal Learning by Arming Image ViTs with Video UniFormer. In ICCV.

- [24] Kunchang Li, Yali Wang, Yizhuo Li, Yi Wang, Yinan He, Limin Wang, and Yu Qiao. 2023. Unmasked Teacher: Towards Training-Efficient Video Foundation Models. In ICCV. 19891–19903.
- [25] Linjie Li, Zhe Gan, Kevin Lin, Chung-Ching Lin, Zicheng Liu, Ce Liu, and Lijuan Wang. 2023. LAVENDER: Unifying Video-Language Understanding as Masked Language Modeling. In CVPR. 23119–23129.
- [26] Yiming Li, Xiangdong Wang, and Hong Liu. 2024. Audio-Free Prompt Tuning for Language-Audio Models. In ICASSP. 491–495.
- [27] Tsung-Yi Lin, Michael Maire, Serge J. Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. 2014. Microsoft COCO: Common Objects in Context. In ECCV, Vol. 8693. 740–755.
- [28] Xiaoyu Liu, Hanlin Lu, Jianbo Yuan, and Xinyu Li. 2023. Cat: causal audio transformer for audio classification. In ICASSP. 1–5.
- [29] Yuzhuo Liu, Xubo Liu, Yan Zhao, Yuanyuan Wang, Rui Xia, Pingchuan Tain, and Yuxuan Wang. 2024. Audio Prompt Tuning for Universal Sound Separation. In ICASSP. 1446–1450.
- [30] Yongcheng Liu, Lu Sheng, Jing Shao, Junjie Yan, Shiming Xiang, and Chunhong Pan. 2018. Multi-Label Image Classification via Knowledge Distillation from Weakly-Supervised Detection. In ACM MM.
- [31] Ziming Liu, Song Guo, Jingcai Guo, Yuanyuan Xu, and Fushuo Huo. 2022. Towards unbiased multi-label zero-shot learning with pyramid and semantic attention. TMM (2022).
- [32] Ziming Liu, Song Guo, Xiaocheng Lu, Jingcai Guo, Jiewei Zhang, Yue Zeng, and Fushuo Huo. 2023. (ML)² P-Encoder: On Exploration of Channel-Class Correlation for Multi-Label Zero-Shot Learning. In CVPR. 23859–23868.
- [33] Ze Liu, Jia Ning, Yue Cao, Yixuan Wei, Zheng Zhang, Stephen Lin, and Han Hu. 2022. Video swin transformer. In CVPR, 3202–3211.
- [34] L. Nanni, Y.M.G. Costa, D.R. Lucio, C.N. Silla, and S. Brahnam. 2017. Combining visual and acoustic features for audio classification tasks. PRL 88 (Mar 2017), 49–56
- [35] Xing Nie, Bolin Ni, Jianlong Chang, Gaofeng Meng, Chunlei Huo, Shiming Xiang, and Qi Tian. 2023. Pro-tuning: Unified Prompt Tuning for Vision Tasks. TCSVT (2023), 1–1.
- [36] Karol J Piczak. 2015. ESC: Dataset for environmental sound classification. In ACM MM. 1015–1018.
- [37] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. 2021. Learning Transferable Visual Models From Natural Language Supervision. In ICML, Vol. 139. 8748–8763.
- [38] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. 2015. Imagenet large scale visual recognition challenge. IJCV 115 (2015), 211–252.
- [39] Justin Salamon, Christopher Jacoby, and Juan Pablo Bello. 2014. A dataset and taxonomy for urban sound research. In ACM MM. 1041–1044.
- [40] Pritam Sarkar and Ali Etemad. 2023. Self-Supervised Audio-Visual Representation Learning with Relaxed Cross-Modal Synchronicity. In AAAI. 9723–9732.
- [41] Linus Scheibenreif, Michael Mommert, and Damian Borth. 2023. Masked vision transformers for hyperspectral image classification. In CVPR. 2165–2175.
- [42] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. 2022. Laion-5b: An open large-scale dataset for training next generation image-text models. NeurIPS 35 (2022), 25278–25294.
- [43] Shuai Shao, Zeming Li, Tianyuan Zhang, Chao Peng, Gang Yu, Xiangyu Zhang, Jing Li, and Jian Sun. 2019. Objects365: A large-scale, high-quality dataset for object detection. In ICCV. 8430–8439.
- [44] Christian Simon, Piotr Koniusz, and Mehrtash Harandi. 2022. Meta-learning for multi-label few-shot classification. In CVPR. 3951–3960.
- [45] Khurram Soomro, Amir Roshan Zamir, and Mubarak Shah. 2012. UCF101: A Dataset of 101 Human Actions Classes From Videos in The Wild. CoRR abs/1212.0402 (2012).
- [46] Ximeng Sun, Ping Hu, and Kate Saenko. 2022. Dualcoop: Fast adaptation to multilabel recognition with limited annotations. In NeurIPS, Vol. 35. 30569–30582.
- [47] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurélien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. LLaMA: Open and Efficient Foundation Language Models. CoRR abs/2302.13971 (2023).
- [48] Du Tran, Heng Wang, Lorenzo Torresani, Jamie Ray, Yann LeCun, and Manohar Paluri. 2018. A Closer Look at Spatiotemporal Convolutions for Action Recognition. In CVPR. 6450–6459.
- [49] Laurens van der Maaten and Geoffrey Hinton. 2008. Visualizing Data using t-SNE. JMLR 9, 86 (2008), 2579–2605.
- [50] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is All you Need. In NeurIPS. 5998–6008.

- [51] Fengqiang Wan, Xiangyu Wu, Zhihao Guan, and Yang Yang. 2024. CoVLR: Coordinating Cross-Modal Consistency and Intra-Modal Relations for Vision-Language Retrieval. In ICME. 1–6.
- [52] Haixin Wang, Jianlong Chang, Yihang Zhai, Xiao Luo, Jinan Sun, Zhouchen Lin, and Qi Tian. 2024. LION: Implicit Vision Prompt Tuning. In AAAI. 5372–5380.
- [53] Junke Wang, Dongdong Chen, Zuxuan Wu, Chong Luo, Luowei Zhou, Yucheng Zhao, Yujia Xie, Ce Liu, Yu-Gang Jiang, and Lu Yuan. 2022. OmniVL: One Foundation Model for Image-Language and Video-Language Tasks. In NeurIPS.
- [54] Meng Wang, Changzhi Luo, Richang Hong, Jinhui Tang, and Jiashi Feng. 2016. Beyond Object Proposals: Random Crop Pooling for Multi-Label Image Recognition. TIP 25, 12 (Dec 2016), 5678–5688.
- [55] Yi Wang, Yinan He, Yizhuo Li, Kunchang Li, Jiashuo Yu, Xin Ma, Xinyuan Chen, Yaohui Wang, Ping Luo, Ziwei Liu, Yali Wang, Limin Wang, and Yu Qiao. 2024. InternVid: A Large-scale Video-Text Dataset for Multimodal Understanding and Generation. In ICLR.
- [56] Yifan Wang, Jie Song, Limin Wang, Luc Van Gool, and Otmar Hilliges. 2016. Two-Stream SR-CNNs for Action Recognition in Videos. In BMCV.
- [57] Junda Wu, Tong Yu, Rui Wang, Zhao Song, Ruiyi Zhang, Handong Zhao, Chaochao Lu, Shuai Li, and Ricardo Henao. 2023. InfoPrompt: Information— Theoretic Soft Prompt Tuning for Natural LanguageUnderstanding. In NeurIPS, Vol. 36.
- [58] Xiangyu Wu, Qing-Yuan Jiang, Yang Yang, Yi-Feng Wu, Qing-Guo Chen, and Jianfeng Lu. 2024. TAI++: Text as Image for Multi-Label Image Classification by Co-Learning Transferable Prompt. In IJCAI.
- [59] Xiangyu Wu, Jianfeng Lu, Zhuanfeng Li, and Fengchao Xiong. 2022. Ques-to-Visual Guided Visual Question Answering. In ICIP. 4193–4197.
- [60] Xiangyu Wu, Feng Yu, Yang Yang, Qing-Guo Chen, and Jianfeng Lu. 2025. Multi-Label Test-Time Adaptation with Bound Entropy Minimization. In The Thirteenth International Conference on Learning Representations, ICLR 2025, Singapore, April 24-28, 2025. OpenReview.net. https://openreview.net/forum?id=75PhjtbBdr
- [61] Yusong Wu, Ke Chen, Tianyu Zhang, Yuchen Hui, Taylor Berg-Kirkpatrick, and Shlomo Dubnov. 2023. Large-Scale Contrastive Language-Audio Pretraining with Feature Fusion and Keyword-to-Caption Augmentation. In ICASSP. 1–5.
- [62] Kele Xu, Kang You, Ming Feng, and Boqing Zhu. 2023. Trust-worth multirepresentation learning for audio classification with uncertainty estimation. JASA 153 (Mar 2023), A125–A125.
- [63] Hongwei Xue, Yuchong Sun, Bei Liu, Jianlong Fu, Ruihua Song, Houqiang Li, and Jiebo Luo. 2022. CLIP-ViP: Adapting Pre-trained Image-Text Model to Video-Language Representation Alignment. CoRR (2022).
- [64] Shen Yan, Xuehan Xiong, Anurag Arnab, Zhichao Lu, Mi Zhang, Chen Sun, and Cordelia Schmid. 2022. Multiview Transformers for Video Recognition. In CVPR. 3323–3333.
- [65] Shen Yan, Xuehan Xiong, Anurag Arnab, Zhichao Lu, Mi Zhang, Chen Sun, and Cordelia Schmid. 2022. Multiview transformers for video recognition. In CVPR. 3333–3343.
- [66] Shuo Yang, Zirui Shang, Yongqi Wang, Derong Deng, Hongwei Chen, Qiyuan Cheng, and Xinxiao Wu. 2024. Data-free Multi-label Image Recognition via LLM-powered Prompt Tuning. CoRR abs/2403.01209 (2024).
- [67] Taojiannan Yang, Yi Zhu, Yusheng Xie, Aston Zhang, Chen Chen, and Mu Li. 2023. AIM: Adapting Image Models for Efficient Video Action Recognition. In ICLR
- [68] Yang Yang, Fengqiang Wan, Qing-Yuan Jiang, and Yi Xu. 2024. Facilitating Multimodal Classification via Dynamically Learning Modality Gap. In NeurIPS.
- [69] Hantao Yao, Rui Zhang, and Changsheng Xu. 2023. Visual-Language Prompt Tuning with Knowledge-Guided Context Optimization. In CVPR. 6757–6767.
- [70] Jiahui Yu, Zirui Wang, Vijay Vasudevan, Legg Yeung, Mojtaba Seyedhosseini, and Yonghui Wu. 2022. CoCa: Contrastive Captioners are Image-Text Foundation Models. TMLR 2022 (2022).
- [71] Xinlei Yu, Changmiao Wang, Hui Jin, Ahmed Elazab, Gangyong Jia, Xiang Wan, Changqing Zou, and Ruiquan Ge. 2025. CRISP-SAM2: SAM2 with Cross-Modal Interaction and Semantic Prompting for Multi-Organ Segmentation. arXiv preprint arXiv:2506.23121 (2025).
- [72] Chunhui Zhang, Xin Sun, Yiqian Yang, Li Liu, Qiong Liu, Xi Zhou, and Yanfeng Wang. 2023. All in One: Exploring Unified Vision-Language Tracking with Multi-Modal Alignment. In ACM MM. 5552-5561.
- [73] Yifei Zhang, Hao Zhu, Zixing Song, Piotr Koniusz, and Irwin King. 2023. Spectral feature augmentation for graph contrastive learning and beyond. In AAAI, Vol. 37. 11289–11297
- [74] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. 2022. Learning to Prompt for Vision-Language Models. IJCV 130, 9 (2022), 2337–2348.
- [75] Bin Zhu, Bin Lin, Munan Ning, Yang Yan, Jiaxi Cui, Hongfa Wang, Yatian Pang, Wenhao Jiang, Junwu Zhang, Zongwei Li, Wancai Zhang, Zhifeng Li, Wei Liu, and Li Yuan. 2024. LanguageBind: Extending Video-Language Pretraining to N-modality by Language-based Semantic Alignment. In ICLR.
- [76] Xuelin Zhu, Jiuxin Cao, Jian Liu, Dongqi Tang, Furong Xu, Weijia Liu, Jiawei Ge, Bo Liu, Qingpei Guo, and Tianyi Zhang. 2023. Text as Image: Learning Transferable Adapter for Multi-Label Classification. CoRR abs/2312.04160 (2023).

[77] Xuelin Zhu, Jian Liu, Weijia Liu, Jiawei Ge, Bo Liu, and Jiuxin Cao. 2023. Scene-aware label graph learning for multi-label image classification. In ICCV. 1473–1402.

A DETAILS OF PRE-TRAINED MULTIMODAL MODELS.

Our TaAM-CPT is built upon multimodal pre-trained models, including video-language model, image-language model, and audio-language model, and uses frozen text encoders for prompt tuning, as well as frozen modality encoders for object recognition predicting. In our work, we choose the pretrained multimodal models, open-sourced by the LAION [42] organization, as the modality-aligned text and modality encoders. For a total of 300k text sentences on a single Tesla V100 for the Kinetic-400, MSCOCO, and ESC50 datasets, each epoch takes 12 minutes and the total training cost for 10 epochs is about 2 hours.

ViCLIP. ViCLIP is a video-language pretraining model, building upon the open-source CLIP of OpenAI. The model consists of a video encoder and corresponding text encoder, which is pretrained on the InternVid dataset containing 7 million videos, each with detailed text descriptions. We use the BASE architecture as our baseline model with 12 attention layers and 512 encoding dimensions

CLIP. We select the open-source image-language pretraining model released by the LAION organization as our baseline model. The model comprises an image encoder and corresponding transformer-based text encoder, each with 12 attention layers and an encoding dimension of 512. The size of the input image is 224×224 , with the patch size being 32. For image modality, CLIP-ViT-B-32 [6] is selected as the image encoder and image-text encoder.

CLAP. For the audio-language pretraining model, likewise, we select CLAP released by the LAION organization as our baseline model. The audio encoder is a transformer-based model with 4 groups of swin-transformer blocks, while the text encoder is RoBERTa. Two-layer MLPs with ReLU activation are applied to mAP both audio and text outputs into 512 dimensions. For audio modality, we select CLAP [61] from LAION [42] as the audio encoder and the built-in Robert as the audio-text encoder.

B DETAILS OF DATASETS

B.1 Video Datasets

UCF101. UCF101 [45] is a commonly used video classification dataset that contains 101 different action classes, each class contains approximately 100~300 video clips, and a total of 13,320 video clips. These video clips are collected from real data on YouTube, ranging in length from 10~30 seconds. We use all of the video data to evaluate our methods.

Kinetic-400. Kinetic-400 [4] is a large-scale, high-quality video dataset collected from YouTube, including 400 human action classes. Each action class contains 450~1150 video clips, covering a wide range of classes, e.g., playing instruments, interactions between humans and objects, and handshakes. Each action has 250~1000 video clips for the training set, 50 video clips for the validation set, and 100 video clips for the test set. The validation set is used to evaluate our methods.

Kinetic-600. Kinetic-600 [2] is an extension of the Kinetic-400 dataset, comprising approximately 480K video clips from 600 action classes. Each action class has at least 700 video clips. The dataset consists of 450~1000 video clips for training, 50 for validation, and

 $100\ {\rm for}\ {\rm testing}\ {\rm per}\ {\rm action}\ {\rm class}.$ The validation set is used to evaluate our methods.

Kinetic-700. Kinetic-700 [3] is an extension of the Kinetic-600 dataset, covering 700 human action classes. Each action class has at least 700 video clips. Each video is a 10-second action clip extracted from original YouTube videos and labeled accordingly. There are a total of 650,000 video clips, with each action class comprising 450,100 video clips for training, 5,000 video clips for validation, and 1,000 video clips for testing. We use the validation set to evaluate our methods.

B.2 Image Datasets

MSCOCO. MSCOCO [27] is a large-scale computer vision dataset used for tasks such as object recognition, object detection, and image segmentation. It includes 80 image classes, 328,000 images, and 2,500,000 instances. It comprises 82,783 training images, 40,504 validation images, and 40,775 test images. We use the validation set to evaluate our methods.

VOC2007. VOC2007 [11] is an image dataset containing 20 image classes that can be used to evaluate image classification, object detection, and image segmentation tasks. It consists of 9,963 images in total, with 5,011 images in the training set and 4,952 images in the test set. The test set is used to evaluate our methods.

VOC2012. VOC2012 [11] dataset contains 20 classes, including people, animals, vehicles, indoor objects, and a background category, making a total of 20 classes. It can be used for evaluating image classification, object detection, and image segmentation tasks. It comprises 11,540 images, with 5,717 images in the training set and 5,823 images in the test set. The test set is used to evaluate our methods.

NUSWIDE. NUSWIDE [7] is an image dataset that contains 269,648 images collected from Flickr, with a total of 81 manually annotated concepts, including objects and scenes. It includes 161,789 images for the training set and 107,859 images for the validation set. We use the validation set to evaluate our methods.

ImageNet-mini. ImageNet-mini [38] is derived from the ImageNet dataset and contains 100 classes with a total of 60,000 images, with 600 samples per class. The training and validation sets are typically divided into an 8:2 ratio by class. (For small sample classification, 64 classes are used for training, 16 for validation, and 20 for testing.) We use the test set to evaluate our methods.

Objects365. Objects365 [43] is a large object detection dataset that contains 638k images, 365 image classes, and 10,101k bounding boxes, far surpassing datasets like COCO. According to the paper's annotation process, a total of 740k images were annotated, with 600k used for training, 38k for validation, and 100k for testing. We use the test set to evaluate our methods.

B.3 Audio Datasets

ESC50. ESC50 [36] is a standard dataset for environmental sound classification that contains 50 different environmental categories, each with 40 samples of up to 5 seconds in duration, totaling 2,000 samples. These samples cover a wide range of environments, such as animal sounds, traffic noise, indoor activities, etc. All samples are carefully balanced to ensure uniformity when training models. We use the validation set to evaluate our methods.

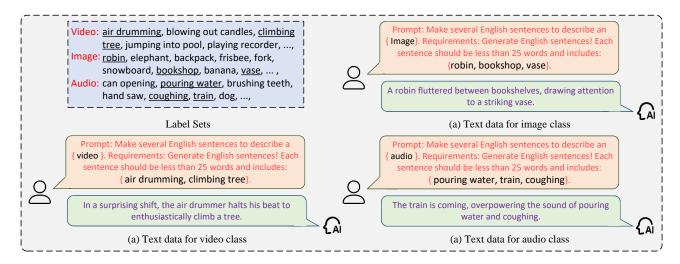


Figure 7: The candidate label set and text data generated by LLMs.

US8K. UrbanSound8k [39] is a widely used open data set for automatic urban environment sound classification, which includes ten categories such as air conditioning sound and car horn sound. There are 8732 audio clips in the dataset with a length of about 4 seconds. The data set is divided into training and testing sets. We use the test set to evaluate our methods.

C TRAINING TEXT DATA CONSTRUCTION.

Here, we discuss the text training data construction for different modalities. We construct the following prompt template to input into LLaMA-2-7B for generating text description data.

TEMPLATE: Make several English sentences to describe a { **Modality** }. Requirements: Generate 5 English sentences! Each sentence should be less than 25 words and includes: { **Labels** }.

where { Modality } is replaced with video, audio, and image, { Labels } denotes the sampled classes. For video modality, which typically involves single classification tasks, we set the number of sampled categories to 2 to prevent too many categories from appearing in one sentence, which could interfere with the model's learning of specific representations for each category. For image classification datasets, where multiple categories can appear on a single image and audio modalities, the number of sampled categories is set to 1, 2, or 3 to ensure that the model not only learns the dependencies between categories but also acquires independent representations for each category. As shown in Figure 7, we randomly select several classes from the label set and construct a prompt template to query the LLMs to generate text data containing the semantic information of these classes.

D FURTHER ANALYSIS

Prompt Design. Here, we mainly discuss the variants of consistent prompt tuning (CPT) in Table 10: a) Shared-Intra (1024), where the prompt is initialized as 1024-d vector and mapped to 512-d through a FC; b) Shared-Intra (512) represents initialization as a 512-d vector and then mapped to 512-d; c) Shared-Inter (512), where all prompts

Table 10: Results of different prompt designs.

Prompt	K400	MSCOCO	ESC50
Shared-Intra (1024)	(43.1, 74.2)	55.4	90.6
Shared-Intra (512)	(47.5, 75.3)	58.7	91.9
Shared-Inter (512)	(50.1, 79.3)	62.2	92.1
TaAM-CPT(Ours)	(55.2, 80.4)	68.1	94.2

across all modalities share a FC and are mapped to 512-d. On Kinetic-400, we note a pronounced degradation of these variants. We believe the decline is mainly attributable to the numerous categories that are semantically proximate (e.g., *making pizza* and *making sandwich*). These phenomena are also observed in the MSCOCO and ESC50 datasets.

Table 11: Different loss weight between intra- and inter-modal learning.

\mathcal{L}_{Ia}	$\mathcal{L}_{\mathrm{Ie}}$	K400	MSCOCO	ESC50
0.4	1.6	(54.9, 80.0)	67.9	94.0
0.8	1.2	(55.1, 80.2)	68.1	94.1
1.0	1.0	(55.2, 80.4)	68.1	94.2
1.2	0.8	(55.0, 80.2)	68.0	94.0
1.6	0.4	(54.5, 79.6)	68.0	93.9

Loss Weight. In this study, we design Ranking loss and unidirectional contrastive loss to perform intra-modal learning and inter-modal learning. The Ranking loss aims to learn class-specific prompt for each modality, while the contrastive loss is applied to guide the learning of weaker modalities (video) through those stronger ones (image and audio). Here, we explore the impact of setting different loss weights for these two loss functions. As shown in Figure 11, \mathcal{L}_{Ia} represents the Ranking loss for intra-modal learning, and \mathcal{L}_{Ie} represents the uni-directional contrastive loss for intermodal learning. Our method achieves the best results when the weights of \mathcal{L}_{Ia} and \mathcal{L}_{Ie} are identical. Additionally, we notice that when the weight of \mathcal{L}_{Ie} is set to 1.0,0.8 and 0.4, there is a significant decrease in top-1 and top-5 accuracy on the Kinetic-400 dataset, while the performance on MSCOCO and ESC50 datasets only suffer

minor damage. This indicates that inter-modal learning greatly affects the learning of weaker modality, which is the video modality in this case.

Table 12: Results of different prompt initialization.

Prompt initialization	K400	MSCOCO	ESC50
ZS-ViCLIP,CLIP,CLAP Initialize by CLIP, w/o $\mathcal{L}_{\mathbf{Inter}}$	(53.8, 78.7) (54.5, 79.6)	55.6 65.3	90.5 93.1
TaAM-CPT(Ours)	(55.2, 80.4)	68.1	94.2

Prompt Initialization. Here, we explore the initializations of the prompt in Table 12. Different from randomly initializing the prompt in the method, we use the output embeddings by CLIP's text encoder to initialize class-specific prompt and remove inter-modal

learning. Therefore, each class-specific prompt encompasses class-specific textual prior knowledge, allowing TaAM-CPT to converge quickly with less training data (we collect only 50 text training data for one class). Although without inter-modal learning, TaAM-CPT achieves higher performance compared to CLIP, ViCLIP, and CLAP.

E VISUALIZATION OF INTRA-MODAL LEARNING.

Here, as shown in Figure 8, 9, 10, 11, 12, we present the more visualization results of the distribution of class-specific prompt learned by intra-modal learning on Kinetic-600/700, MSCOCO, ImageNet-mini, and ESC50 datasets.

Received 20 February 2007; revised 12 March 2009; accepted 5 June 2009

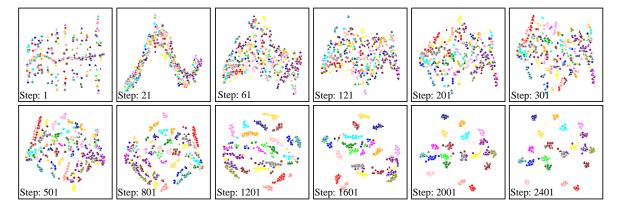


Figure 8: Visualization of the distribution of video prompt and video feature using t-SNE [49] for dimensionality reduction. We randomly select 20 video classes from the Kinetic-600 dataset.

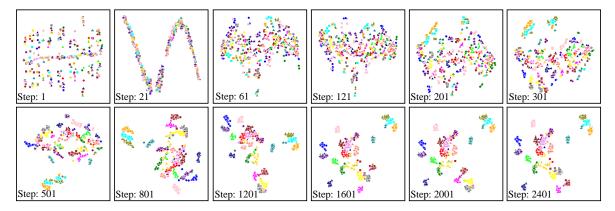


Figure 9: Visualization of the distribution of video prompt and video feature using t-SNE [49] for dimensionality reduction. We randomly select 20 video classes from the Kinetic-700 dataset.

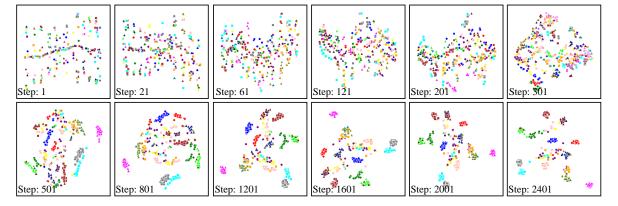


Figure 10: Visualization of the distribution of image prompt and image feature using t-SNE [49] for dimensionality reduction. We randomly select 20 image classes from the MSCOCO dataset.

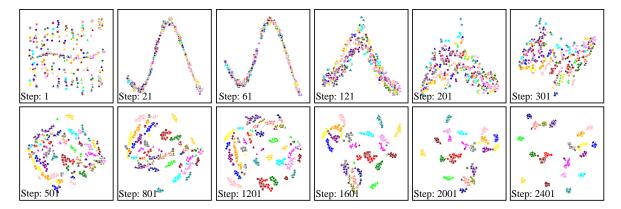


Figure 11: Visualization of the distribution of image prompt and image feature using t-SNE [49] for dimensionality reduction. We randomly select 20 image classes from the ImageNet-mini dataset.

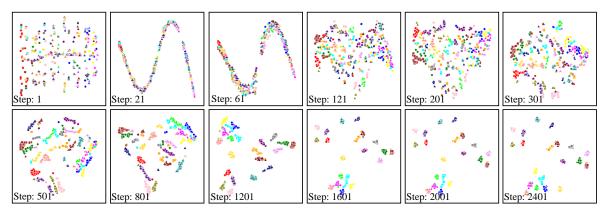


Figure 12: Visualization of the distribution of audio prompt and audio feature using t-SNE [49] for dimensionality reduction. We randomly select 20 audio classes from the ESC50 dataset.