

Bridging Modalities via Progressive Re-alignment for Multimodal Test-Time Adaptation

Jiacheng Li^{1, 3}, Songhe Feng^{2, 3*}

¹Key Laboratory of Big Data & Artificial Intelligence in Transportation (Beijing Jiaotong University),
Ministry of Education, China

²Tangshan Research Institute, Beijing Jiaotong University, China

³School of Computer Science and Technology, Beijing Jiaotong University, Beijing, China
jiacheng.li@bjtu.edu.cn, shfeng@bjtu.edu.cn

Abstract

Test-time adaptation (TTA) enables online model adaptation using only unlabeled test data, aiming to bridge the gap between source and target distributions. However, in multimodal scenarios, varying degrees of distribution shift across different modalities give rise to a complex coupling effect of unimodal shallow feature shift and cross-modal high-level semantic misalignment, posing a major obstacle to extending existing TTA methods to the multimodal field. To address this challenge, we propose a novel multimodal test-time adaptation (MMTTA) framework, termed as **Bridging Modalities via Progressive Re-alignment (BriMPR)**. BriMPR, consisting of two progressively enhanced modules, tackles the coupling effect with a divide-and-conquer strategy. Specifically, we first decompose MMTTA into multiple unimodal feature alignment sub-problems. By leveraging the strong function approximation ability of prompt tuning, we calibrate the unimodal global feature distributions to their respective source distributions, so as to achieve the initial semantic re-alignment across modalities. Subsequently, we assign the credible pseudo-labels to combinations of masked and complete modalities, and introduce inter-modal instance-wise contrastive learning to further enhance the information interaction among modalities and refine the alignment. Extensive experiments on MMTTA tasks, including both corruption-based and real-world domain shift benchmarks, demonstrate the superiority of our method.

Code — <https://github.com/Luchicken/BriMPR>

Introduction

Despite the remarkable success of deep neural networks in various fields, their excellent performances often hinge on specific data conditions. The possible distribution shift (or domain shift) between training and testing data has become a major obstacle to model generalization. Unsupervised domain adaptation (UDA) (Tzeng et al. 2014; Long et al. 2015; Jin et al. 2020) and domain generalization (DG) (Zhou et al. 2021; Li et al. 2018; Zhang and Feng 2023) have been proposed to mitigate domain gaps by designing sophisticated strategies that enable the model to adapt to the target domain during training. In contrast, test-time adaptation (TTA) (Sun

et al. 2020; Wang et al. 2021; Zhang, Levine, and Finn 2022; Niu et al. 2022) adjusts the model according to specific test data during the test stage, reducing the dependence on the training process and training data, thereby making it a promising and more practical solution.

With the advancement of sensor technology, integrating and leveraging multimodal data collected from diverse sensors has significantly enhanced the perception capability of intelligent systems. Nevertheless, multimodal data also suffer from distribution shifts. What’s worse, due to the complexity of multimodal data, different modalities often exhibit varying degrees of distribution shift from the source domain, inducing a complex coupling effect of unimodal shallow feature shift and cross-modal high-level semantic misalignment. Existing TTA methods, which are primarily designed for unimodal tasks, struggle to ensure consistent improvements across all modalities and often fail to fully exploit the rich information available in multimodal inputs. In Fig. 1, we visualize both unimodal and multimodal feature representations during the adaptation on the audio-visual event classification dataset Kinetics50-C (Yang et al. 2024). As a representative unimodal TTA method, EATA (Niu et al. 2022) reduces the uncertainty of model predictions by minimizing the entropy of reliable samples. However, it shows limited improvement in bridging the domain gap between source and target features for each modality. READ (Yang et al. 2024), a pioneering method for multimodal test-time adaptation (MMTTA), adapts the model by updating the self-attention layers in the fusion module to assign more weights to the high-quality modality. Nevertheless, it lacks the correction of shallow unimodal features. As shown in Fig. 1a and Fig. 1b, the lack of effective guidance for unimodal features hinders proper alignment across modalities. As a result, the fused multimodal feature representations derived from multiple unimodal features become entangled, leading to a significant decline in discriminability.

In this work, we propose **Bridging Modalities via Progressive Re-alignment (BriMPR)** for multimodal test-time adaptation. Through the joint efforts of self-calibration for each modality and inter-modal information interaction, BriMPR realigns the modalities that are subject to distribution shift with each other. Since the feature representations of each modality are well-aligned in the source space, we first decompose MMTTA into multiple unimodal feature

*Corresponding Author.

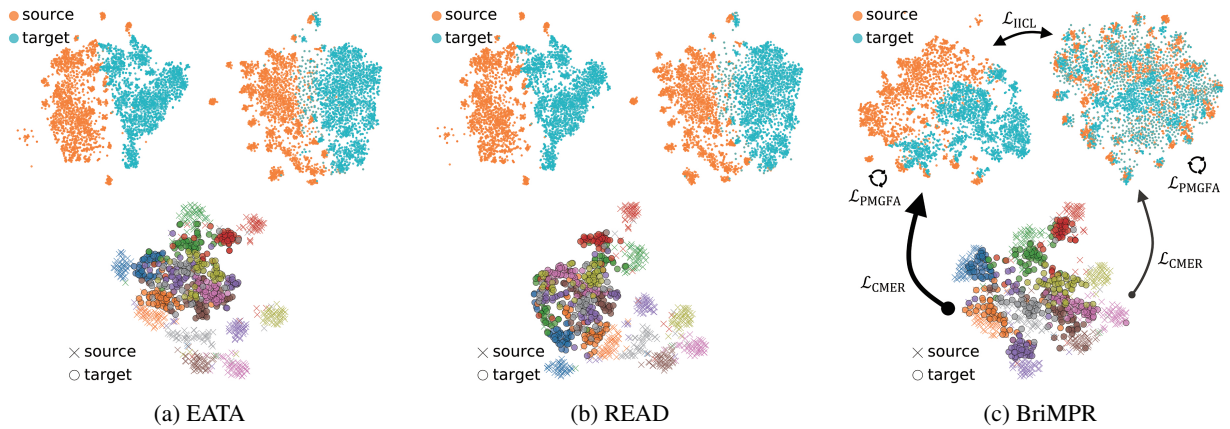


Figure 1: t-SNE visualizations of unimodal (top) and fused multimodal (bottom) features during adaptation versus source features. For fused features, 10 classes from Kinetics50-C are shown.

alignment sub-problems. Leveraging the strong function approximation ability of prompt tuning (Wang et al. 2023), we calibrate the global feature distribution of each modality to its corresponding source distribution via modality-specific prompts embedded across layers of the modality-specific encoders, thereby indirectly achieving initial cross-modal semantic alignment. Subsequently, the alignment is further refined by enhancing inter-modal information interaction. We propose a novel cross-modal masked embedding recombination loss, which promotes the extraction of multimodal information by providing calibrated pseudo-labels for the combinations of masked and complete modalities. Additionally, we introduce inter-modal instance-wise contrastive learning to maintain cross-modal alignment at the instance level. As shown in Fig. 1c, BriMPR effectively bridges the domain gap between the source and target for each unimodal feature, thereby enhancing the discriminability of the fused features. Our contributions can be summarized as follows:

- We propose a novel MMTA framework which mitigates modality-wise distribution shifts in a divide-and-conquer manner, facilitating the re-alignment among modalities.
- We leverage the excellent function approximation ability of prompt tuning to achieve efficient calibration of the unimodal global feature distribution, and propose a novel cross-modal masked embedding recombination strategy to enhance the inter-modal interaction.
- We conduct extensive experiments on MMTA benchmarks, including corruption shift and real-world shift datasets, demonstrating the superiority of BriMPR over existing SOTA methods.

Related Work

Test-Time Adaptation. Test-time adaptation (TTA) leverages unlabeled test data to adapt models to unseen target domains during test-time. The idea of TTA can be traced back to TTT (Sun et al. 2020), which uses a self-supervised auxiliary branch to enable adaptation during inference. A series of works (Wang et al. 2021; Niu et al. 2022, 2023; Lee et al. 2024) explore fully test-time adaptation (FTTA)

by optimizing the normalization layers via entropy-based losses, without altering the pre-training stage. Given the limitations of unimodal TTA methods in multimodal scenarios, MM-TTA (Shin et al. 2022) proposes a cross-modal self-learning framework for MMTA. READ (Yang et al. 2024) highlights the reliability bias of MMTA under unimodal corruption, and proposes to adaptively assign modality weights by optimizing the self-attention in the fusion module. ABPEM (Zhao et al. 2025) reduces the gap between cross-attention and self-attention, and computes the principal part of entropy to reduce gradient noise. SuMi (Guo and Jin 2025) utilizes interquartile range smoothing to identify samples used for calculating entropy loss. Moreover, AEO (Dong, Chatzi, and Fink 2025) introduces unseen classes and proposes the Multimodal open-set test-time adaptation setting. In this work, we attribute the difficulties of MMTA to the coupling effect of unimodal shallow feature shift and cross-modal high-level semantic misalignment, and propose a divide-and-conquer method to re-bridge modalities during testing.

Prompt Tuning. Originally developed in natural language processing, prompt tuning introduces extra tokens to guide models toward generating task-specific outputs. In computer vision, approaches like CoOp (Zhou et al. 2022b) and Co-CoOp (Zhou et al. 2022a) leverage learnable prompts to enhance the zero-shot recognition capabilities of vision-language models (VLMs). Integrating the idea of TTA, test-time prompt tuning (TPT) (Shu et al. 2022; Feng et al. 2023; Zhang et al. 2024a) fine-tunes text prompts using test samples to improve the generalization of VLMs. While TPT primarily focuses on extracting rich knowledge from large-scale VLMs, our work is more closely aligned with visual prompt tuning (VPT) (Jia et al. 2022; Yoo et al. 2023). VPT introduces prompt tuning into Vision Transformer, achieving significant performance gains over full fine-tuning. Our work extends prompt tuning to MMTA tasks, leveraging the strong function approximation ability of prompts to efficiently calibrate the distribution of each unimodal feature—not limited to visual features alone.

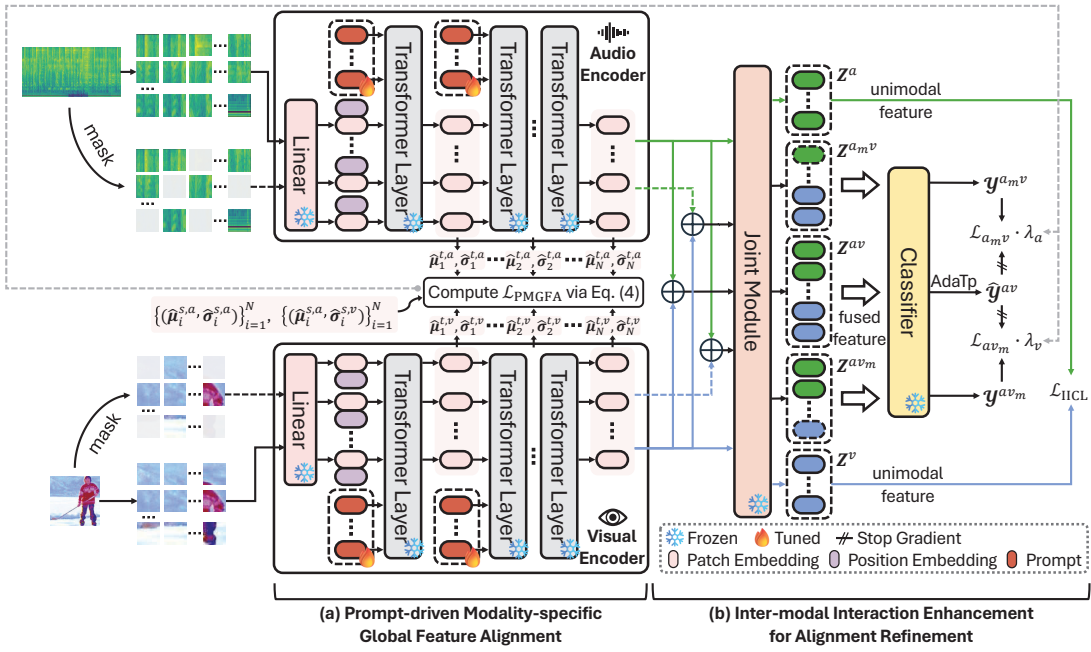


Figure 2: Overview of BriMPR. BriMPR achieves initial alignment and alignment refinement through two progressive modules. The added modality-specific prompts are used to project the unimodal features into the re-aligned feature space.

Preliminaries

Multimodal Test-Time Adaptation (MMTTA). Without loss of generality, we take two modalities as an example to provide a formal definition of MMTTA. An off-the-shelf model \mathcal{F}_Θ pre-trained on the source domain $\mathcal{D}_S = \{(\mathbf{x}_i^{u_1}, \mathbf{x}_i^{u_2}, y_i)\}_{i=1}^{N_S}$ is adopted as the initial model, where the two modalities of the source data follow the probability distributions $\mathbf{x}_i^{u_1} \sim P_{S,u_1}(\mathbf{x})$ and $\mathbf{x}_i^{u_2} \sim P_{S,u_2}(\mathbf{x})$, respectively. The goal of MMTTA is to adapt \mathcal{F}_Θ online to the target domain $\mathcal{D}_T = \{(\mathbf{x}_j^{u_1}, \mathbf{x}_j^{u_2})\}_{j=1}^{N_T}$, where the two modalities of target data follow the probability distributions $\mathbf{x}_j^{u_1} \sim P_{T,u_1}(\mathbf{x})$ and $\mathbf{x}_j^{u_2} \sim P_{T,u_2}(\mathbf{x})$. During adaptation, the source domain is inaccessible and there is a domain shift between the source and target distributions, i.e., $P_{S,u_1}(\mathbf{x}) \neq P_{T,u_1}(\mathbf{x})$ and $P_{S,u_2}(\mathbf{x}) \neq P_{T,u_2}(\mathbf{x})$.

Prompt Tuning. Prompt tuning is regarded as a parameter-efficient fine-tuning technique, which adapts the model to downstream tasks by prepending and optimizing learnable prompt tokens into the input sequence (Li and Liang 2021; Lester, Al-Rfou, and Constant 2021; Jia et al. 2022; Liu et al. 2022). For an encoder Φ consisting of N transformer layers, when inserting a specified number of prompts into the input sequence at each layer, the forward process of the i -th layer can be formulated as:

$$[-; \mathbf{E}_i] = L_i([-; \mathbf{P}_{i-1}; \mathbf{E}_{i-1}]), \quad i = 1, \dots, N. \quad (1)$$

Here $\mathbf{E}_i = [e_{i,1}; e_{i,2}; \dots; e_{i,m}]$ and $\mathbf{P}_i = [p_{i,1}; p_{i,2}; \dots; p_{i,m_p}]$ denote the sequences of original input tokens and inserted prompt tokens, where m and m_p is the number of tokens, and the token dimension is d . $[-; \cdot]$ denotes token-level concatenation. Then, a supervised loss \mathcal{L} is minimized over

the downstream dataset \mathcal{D}_{ds} to obtain the optimal prompt $\mathbf{P}^* = \{\mathbf{P}_0^*, \mathbf{P}_1^*, \dots, \mathbf{P}_{N-1}^*\}$:

$$\mathbf{P}^* = \arg \min_{\mathbf{P}} \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}_{ds}} \mathcal{L}(h(\text{MeanPool}(\mathbf{E}_N)), y), \quad (2)$$

where h denotes the classifier. In MMTTA, due to the absence of annotation for the test data, the loss must be reformulated to enable the learning of task-specific prompts.

Methodology

In this section, we introduce BriMPR for MMTTA, with its overall framework illustrated in Fig. 2. BriMPR comprises two progressively enhanced modules: (a) *Prompt-driven Modality-specific Global Feature Alignment* achieves initial cross-modal alignment by minimizing the discrepancy between the unimodal target statistics and their corresponding in-distribution statistics; (b) *Inter-modal Interaction Enhancement for Alignment Refinement* further refines the alignment by providing credible pseudo-labels for combinations of masked and complete modalities, and conducting inter-modal instance-wise contrastive learning.

Following READ (Yang et al. 2024), we decompose the source model into two modality-specific encoders (Φ^a for the audio modality and Φ^v for the visual modality), a joint module Ψ , and a classifier h . We update only the prompts for each modality-specific encoder, keeping the rest of the model frozen, to recalibrate individual feature distributions and achieve bottom-up modality re-alignment.

Prompt-driven Modality-specific Global Feature Alignment (PMGFA)

The final prediction of a multimodal model comes from the joint effect of multiple individual modalities. This naturally

allows MMTTA to be decomposed into multiple unimodal test-time adaptation problems. On the other hand, if the target representations at test time can be well projected back to the corresponding source representations, then a TTA model tends to perform well. Based on the intuitions above, we decouple MMTTA into multiple modality-specific feature alignment sub-problems. Since the inter-modal semantic representations are well aligned in the source representation space, solving these sub-problems means indirectly achieving cross-modal semantic alignment of the target representation.

Concretely, we first model the modality-specific source and target feature distributions as multivariate Gaussian distributions, i.e., $P_{S,u} = \mathcal{N}(\mu^{s,u}, \Sigma^{s,u})$ and $P_{T,u} = \mathcal{N}(\mu^{t,u}, \Sigma^{t,u})$, where $u \in \{a, v\}$. In prior works (Liu et al. 2021; Su, Xu, and Jia 2022; Zhang et al. 2024b), feature alignment is typically achieved by matching the first and second moments between distributions (i.e., $\|\mu^t - \mu^s\|_2^2 + \|\Sigma^t - \Sigma^s\|_F^2$) or minimizing the KL-divergence (i.e., $D_{KL}(P_S||P_T)$). However, both approaches rely on the estimation of the covariance matrix Σ , whose error is significantly amplified in high-dimensional data. Therefore, we propose to retain only the diagonal elements of Σ , which reduces the estimation error by a factor of d , as supported by the following theorem:

Theorem 1. *Given $x_1, \dots, x_n \in \mathbb{R}^d$ independently drawn from a multivariate normal distribution $\mathcal{N}(\mu, \Sigma)$, let $\hat{\Sigma}$ be the unbiased sample covariance matrix and $\hat{\sigma}^2 = [\hat{\sigma}_1^2, \dots, \hat{\sigma}_d^2]^T$ be the vector of its diagonal entries. Then, the mean squared errors satisfy:*

$$\mathbb{E} [\|\hat{\Sigma} - \Sigma\|_F^2] = \mathcal{O}\left(\frac{d^2}{n}\right), \mathbb{E} [\|\hat{\sigma}^2 - \sigma^2\|_2^2] = \mathcal{O}\left(\frac{d}{n}\right). \quad (3)$$

Due to space limitations, the corresponding proof can be found in Appendix. Emerging research (Wang et al. 2023) has shown that prompt tuning can serve as universal approximators for sequence-to-sequence functions. Motivated by this, we employ prompts as an implicit mapping from the target feature space to the source feature space. For the data \mathbf{x}^u and the i -th layer of the modality-specific encoder Φ^u , the input sequence $\mathbf{E}_{i-1}^u(\mathbf{x}^u)$ undergoes attention interaction with the added prompts \mathbf{P}_{i-1}^u to obtain the transformed output sequence $\mathbf{E}_i^u(\mathbf{x}^u)$. The global feature representation can be expressed as $\mathbf{Z}_i^u(\mathbf{x}^u) = \text{MeanPool}(\mathbf{E}_i^u(\mathbf{x}^u))$. Subsequently, we minimize the following empirical risk on the current batch $\{(\mathbf{x}_j^a, \mathbf{x}_j^v)\}_{j=1}^B$:

$$\begin{aligned} \mathcal{L}_{\text{PMGFA}} &= \sum_{u \in \{a, v\}} \text{Disc}(P_{S,u}, P_{T,u}) \\ &= \sum_{u \in \{a, v\}} \frac{1}{N} \sum_{i=1}^N (\|\hat{\mu}_i^{t,u} - \hat{\mu}_i^{s,u}\|_2 + \|\hat{\sigma}_i^{t,u} - \hat{\sigma}_i^{s,u}\|_2), \end{aligned} \quad (4)$$

where $\text{Disc}(\cdot, \cdot)$ denotes the mean of the layer-wise distribution discrepancy. For convenience, we will interchangeably use Disc^u and $\text{Disc}(P_{S,u}, P_{T,u})$ in the following context. $\|\cdot\|_2$ denotes the Euclidean norm. $\hat{\mu}_i^{t,u} = \sum_{j=1}^B \mathbf{Z}_i^u(\mathbf{x}_j^u)/B$

and $\hat{\sigma}_i^{t,u} = \sqrt{\sum_{j=1}^B [(\mathbf{Z}_i^u(\mathbf{x}_j^u) - \hat{\mu}_i^{t,u})^2]/(B-1)}$ are the estimated mean and standard deviation, respectively. Similar to many other TTA methods (Niu et al. 2022; Döbler, Marsden, and Yang 2023; Wang et al. 2025), we pre-compute $\{\hat{\mu}_i^{s,u}, \hat{\sigma}_i^{s,u}\}_{i=1}^N$ offline prior to the test phase, and this process is performed only once.

Inter-modal Interaction Enhancement for Alignment Refinement

After initial cross-modal semantic alignment via unimodal feature calibration, we further improve the quality of alignment by inter-modal interactions. By recombining masked and complete modalities, the unmasked low-quality modality is forced to draw multimodal information from credible pseudo-labels. Meanwhile, inter-modal instance-wise contrastive learning is applied to strengthen the alignment across instances.

Cross-modal Masked Embedding Recombination.

Masked language modeling (Devlin et al. 2019) and masked image modeling (He et al. 2022) force model to reconstruct the masked regions by utilizing contextual clues and have been widely used as powerful self-supervised learning paradigms in natural language processing and computer vision tasks, respectively. Related but distinct, our proposed Cross-modal Masked Embedding Recombination (CMER) uses masking to simulate distribution shifts from missing patches, serving as a form of data augmentation.

For input \mathbf{x}^u , we randomly mask a portion (e.g., 50%) of its patches and encode the unmasked part \mathbf{x}^{u_m} using Φ^u with modality-specific prompts \mathbf{P}^u to obtain the masked embedding $\Phi^u(\mathbf{x}^{u_m})$. Then, $\Phi^u(\mathbf{x}^{u_m})$ is recombined with complete embeddings from other modalities and passed to the joint module, generating an augmented representation that simulates unimodal corruption. Taking the masked audio modality as an example, the recombined representations and their predictions are formulated as:

$$\begin{aligned} \mathbf{Z}^{a_m v} &= \Psi([\Phi^a(\mathbf{x}^{a_m}); \Phi^v(\mathbf{x}^v)]), \\ \mathbf{y}^{a_m v} &= \sigma(h(\text{MealPool}(\mathbf{Z}^{a_m v}))), \end{aligned} \quad (5)$$

where σ denotes the softmax function. With the initial alignment from PMGFA, we can utilize the complete multimodal data to provide reliable pseudo-labels for augmented inputs. As pseudo-labels become more reliable in the later stages of adaptation, we further calibrate them via temperature scaling (Hinton, Vinyals, and Dean 2015; Guo et al. 2017):

$$\hat{\mathbf{y}}_k^{av} = \frac{\exp([h(\text{MealPool}(\mathbf{Z}^{av}))]_k / \text{AdaTp})}{\sum_{k'=1}^C \exp([h(\text{MealPool}(\mathbf{Z}^{av}))]_{k'} / \text{AdaTp})}. \quad (6)$$

Here, k and k' denote the k -th and k' -th elements of the tensor, and C represents the number of classes. $\text{AdaTp} = 1 + \tau_0 / (1 + \exp(D_0 - \text{Disc}^J)) \in (1, 1 + \tau_0)$ is the adaptive temperature coefficient, where Disc^J is the distribution discrepancy calculated for the joint module, and τ_0 and D_0 are predefined hyperparameters. When Disc^J is large, AdaTp approaches $1 + \tau_0$ to alleviate overconfident predictions. As Disc^J decreases, AdaTp approaches 1, and Eq. (6) approximates the vanilla softmax function. Subsequently, minimize

Method	Noise			Blur				Weather				Digital				Avg.
	Gauss.	Shot	Impul.	Defoc.	Glass	Motion	Zoom	Snow	Frost	Fog	Bright.	Contr.	Elast.	Pixel.	Jpeg	
Source	48.2	50.0	49.2	67.7	61.6	70.6	66.1	60.9	60.7	44.7	75.9	51.8	65.5	68.7	66.1	60.5
• Tent _{ICLR2021}	48.2	49.8	48.7	67.7	62.1	70.8	67.2	61.8	61.4	33.7	76.0	51.2	66.6	69.6	66.9	60.1
• EATA _{ICML2022}	48.7	50.4	49.6	67.8	63.2	70.8	67.5	62.5	62.5	47.9	76.1	52.2	66.9	69.7	67.4	61.5
• SAR _{ICLR2023}	48.5	50.2	49.2	67.8	63.8	70.9	67.9	63.1	62.7	38.7	76.1	52.2	67.1	69.8	67.4	61.0
• DeYO _{ICLR2024}	48.6	50.2	49.4	67.9	62.6	70.9	67.4	62.5	62.3	40.4	76.1	52.2	66.8	69.8	67.3	61.0
• FOA _{ICML2024}	49.2	50.8	49.7	66.0	65.5	69.8	67.4	62.8	65.7	60.3	74.9	51.9	69.5	68.8	68.0	62.7
• READ [†] _{ICLR2024}	50.7	52.2	51.4	67.9	65.3	71.1	68.7	64.0	65.8	56.3	76.3	53.6	68.7	70.0	68.6	63.4
• ABPEM [†] _{AAAI2025}	<u>52.1</u>	<u>53.1</u>	<u>52.8</u>	69.0	<u>65.6</u>	<u>71.8</u>	<u>68.8</u>	64.1	65.7	57.9	<u>76.6</u>	54.3	69.2	71.1	<u>69.2</u>	<u>64.1</u>
• SuMi [†] _{ICLR2025}	50.1	50.7	50.4	<u>68.2</u>	<u>65.6</u>	72.2	69.7	<u>65.7</u>	67.0	56.5	77.1	<u>55.2</u>	69.3	<u>71.2</u>	68.9	63.9
• BriMPR [†]	55.3	56.1	56.7	67.8	67.9	70.6	<u>68.8</u>	65.9	<u>66.2</u>	64.1	76.2	56.3	72.0	73.7	70.5	65.9
Source	52.9	53.0	53.1	57.2	57.2	58.5	57.5	56.5	57.1	55.6	59.2	53.7	57.1	56.4	57.3	56.2
• Tent _{ICLR2021}	53.2	53.3	53.3	56.8	56.6	57.9	57.2	55.9	56.6	56.5	58.5	53.9	57.5	56.8	56.9	56.1
• EATA _{ICML2022}	53.4	53.5	53.5	57.0	57.0	58.3	57.7	56.3	57.0	56.8	59.1	54.2	57.9	57.2	57.2	56.4
• SAR _{ICLR2023}	53.3	53.3	53.3	56.4	56.5	57.9	57.3	55.6	56.4	56.3	58.8	53.7	57.8	56.9	57.0	56.0
• DeYO _{ICLR2024}	53.3	53.4	53.4	56.7	56.7	58.0	57.3	56.0	56.8	56.4	58.7	53.9	57.7	57.0	57.0	56.2
• FOA _{ICML2024}	52.7	52.7	52.7	53.2	53.6	53.6	53.8	53.4	53.4	53.3	55.6	52.5	55.3	53.7	54.4	53.6
• READ [†] _{ICLR2024}	53.8	54.0	53.8	58.0	57.9	59.2	58.7	57.1	58.2	50.0	60.0	55.2	58.5	<u>57.7</u>	58.2	56.7
• ABPEM [†] _{AAAI2025}	46.5	46.7	46.5	54.2	55.1	56.4	55.2	51.3	53.2	52.1	56.6	52.1	54.4	51.7	54.7	52.4
• SuMi [†] _{ICLR2025}	<u>54.0</u>	<u>54.3</u>	<u>53.8</u>	58.2	<u>58.4</u>	59.4	<u>58.7</u>	<u>57.5</u>	58.2	<u>57.6</u>	59.4	<u>54.8</u>	<u>59.0</u>	57.5	<u>58.2</u>	<u>57.3</u>
• BriMPR [†]	54.9	55.0	55.0	57.9	58.5	58.9	58.7	57.5	<u>58.0</u>	58.5	60.3	54.5	59.7	59.3	59.0	57.7

Table 1: Comparison with SOTA methods on Kinetics50-C (top) and VGGSound-C (bottom) under the unimodal shift setting (severity level 5 of video corruption). [†]Multimodal test-time adaptation methods.

the cross-entropy between the calibrated pseudo-label and the augmented predictions:

$$\begin{aligned}\mathcal{L}_{\text{CMER}} &= \lambda^a \mathcal{L}_{a_m v} + \lambda^v \mathcal{L}_{a v_m} \\ &= -\lambda^a \sum_{k=1}^C \hat{\mathbf{y}}_k^{av} \log \mathbf{y}_k^{a_m v} - \lambda^v \sum_{k=1}^C \hat{\mathbf{y}}_k^{av} \log \mathbf{y}_k^{a v_m},\end{aligned}\quad (7)$$

where $\lambda^u = 1 - \text{Disc}^u / (\text{Disc}^a + \text{Disc}^v)$ ($u \in \{a, v\}$) is the weight of the corresponding term, assigning a higher weight to the augmentation with a milder distribution shift in the masked modality. Intuitively, $\mathcal{L}_{\text{CMER}}$ deliberately discards high-quality modality information, forcing the corrupted modality to independently derive the correct result.

Inter-modal Instance-wise Contrastive Learning. Contrastive learning (He et al. 2020; Chen et al. 2020b) has emerged as a key paradigm in cross-modal representation learning, aiming to improve the quality of representations by aligning the feature spaces of the same semantic instance across different modalities/views. Building upon the calibration of unimodal feature distributions, BriMPR introduces inter-modal instance-wise contrastive learning. For data \mathbf{x}^u ($u \in \{a, v\}$), its unimodal representation is as follows:

$$\mathbf{Z}^u = \Psi(\Phi^u(\mathbf{x}^u)). \quad (8)$$

Subsequently, different unimodal representations of the same instance are regarded as positive pairs, while the others as negative pairs. The contrastive loss is defined as:

$$\mathcal{L}_{\text{IICL}} = -\frac{1}{2B} \sum_{j=1}^B \sum_{u_1 \neq u_2} \log \frac{e^{\text{sim}(\mathbf{Z}_j^{u_1}, \mathbf{Z}_j^{u_2})/\tau}}{\sum_{j'=1}^B e^{\text{sim}(\mathbf{Z}_j^{u_1}, \mathbf{Z}_{j'}^{u_2})/\tau}}, \quad (9)$$

where $\text{sim}(\cdot, \cdot)$ denotes the cosine similarity function, and τ denotes the temperature hyperparameter.

Overall Procedure

To brief, BriMPR optimizes the added modality-specific prompts by minimizing the following loss:

$$\mathcal{L}_{\text{BriMPR}} = \mathcal{L}_{\text{PMGFA}} + \mathcal{L}_{\text{CMER}} + \mathcal{L}_{\text{IICL}}. \quad (10)$$

Experiments

Experimental Setups

Datasets and models. We evaluate our method on four commonly used multimodal datasets, including Kinetics50-C, VGGSound-C (Yang et al. 2024), CMU-MOSI (Zadeh et al. 2016), and CH-SIMS (Yu et al. 2020). Kinetics50-C/VGGSound-C contain two modalities: video and audio, and are obtained by adding various corruptions to the test sets of the original versions (i.e., Kinetics (Kay et al. 2017) and VGGSound (Chen et al. 2020a)). For the video modality and the audio modality, 15 and 6 types of corruption are introduced, respectively, which are divided into 5 severity levels. Following (Yang et al. 2024), we use the pre-trained CAV-MAE (Gong et al. 2023) as the source model. CMU-MOSI/CH-SIMS contain three modalities: text, video, and audio. Following (Guo and Jin 2025), we use stacked Transformer blocks as the backbone and pre-train the model on MOSI and SIMS, respectively.

Considered settings. For domain shifts caused by corruptions, we consider two tasks and report average classification accuracy (%): (1) Under the unimodal shift setting, following (Yang et al. 2024), one modality is corrupted while the other modality remains clean; (2) Under the multimodal shift setting, both modalities are corrupted. For real-world

Method	Noise			Weather			Avg.	Noise			Weather			Avg.
	Gauss.	Traff.	Crowd	Rain	Thund.	Wind		Gauss.	Traff.	Crowd	Rain	Thund.	Wind	
Source	74.3	65.3	68.0	70.3	68.0	70.5	69.4	37.3	21.2	16.9	21.8	27.3	25.7	25.0
• Tent _{ICLR2021}	74.6	67.4	69.5	70.8	67.6	71.2	70.2	10.8	2.8	1.8	2.9	5.6	3.9	4.6
• EATA _{ICML2022}	74.6	67.3	69.4	70.8	69.8	71.0	70.5	40.2	30.0	27.8	29.7	36.5	32.2	32.7
• SAR _{ICLR2023}	74.6	67.0	69.2	70.9	69.5	70.9	70.3	30.4	5.5	8.0	9.3	32.5	17.2	17.1
• DeYO _{ICLR2024}	74.6	67.0	69.3	70.8	69.0	71.0	70.3	22.9	4.9	15.8	4.9	16.5	20.0	14.2
• FOA _{ICML2024}	73.8	70.0	70.5	71.0	73.0	71.2	71.6	31.5	26.2	23.7	31.0	34.2	26.7	28.9
• READ [†] _{ICLR2024}	74.8	69.2	69.9	71.4	72.4	71.0	71.5	39.9	29.4	26.8	30.8	36.8	30.7	32.4
• ABPEM [†] _{AAAI2025}	74.7	68.5	70.3	71.7	72.3	71.2	71.4	38.5	27.6	25.2	26.5	32.7	26.5	29.5
• SuMi [†] _{ICLR2025}	75.1	68.9	<u>70.6</u>	<u>71.6</u>	<u>72.8</u>	72.1	<u>71.9</u>	41.9	26.3	<u>27.9</u>	<u>31.6</u>	<u>37.1</u>	<u>34.1</u>	<u>33.2</u>
• BriMPR [†]	<u>74.8</u>	<u>69.6</u>	71.7	71.5	72.4	<u>72.0</u>	72.0	39.3	35.0	36.7	32.5	41.0	34.6	36.5

Table 2: Comparison with SOTA methods on Kinetics50-C (left) and VGGSound-C (right) under the unimodal shift setting (severity level 5 of audio corruption).

Method	Noise			Blur				Weather				Digital				Avg.
	Gauss.	Shot	Impul.	Defoc.	Glass	Motion	Zoom	Snow	Frost	Fog	Bright.	Contr.	Elast.	Pixel.	Jpeg	
Source	13.1	14.1	13.3	37.2	37.4	45.3	41.8	29.4	32.6	20.4	55.2	18.3	42.5	38.8	37.8	31.8
• Tent _{ICLR2021}	9.1	9.7	9.1	32.5	34.1	43.5	40.2	23.2	28.3	13.2	55.1	13.7	40.7	34.7	35.0	28.1
• EATA _{ICML2022}	12.9	14.0	13.1	38.1	38.7	46.9	43.1	30.6	33.0	20.2	56.5	18.2	43.7	40.7	39.0	32.6
• SAR _{ICLR2023}	11.8	12.8	11.9	37.4	38.2	46.3	43.1	29.8	33.0	17.4	56.0	16.0	43.5	39.5	38.2	31.7
• DeYO _{ICLR2024}	11.0	12.0	11.1	37.0	37.7	46.3	43.2	29.9	33.3	17.9	56.2	17.0	43.7	39.7	38.0	31.6
• FOA _{ICML2024}	18.7	20.6	19.3	43.7	45.5	50.2	<u>47.9</u>	38.9	43.7	37.2	60.5	23.5	<u>52.7</u>	<u>48.9</u>	<u>47.4</u>	<u>39.9</u>
• READ [†] _{ICLR2024}	14.5	14.9	14.8	<u>43.8</u>	42.1	<u>51.0</u>	46.5	35.4	38.9	27.6	58.9	22.6	47.1	42.1	38.1	35.9
• ABPEM [†] _{AAAI2025}	<u>19.2</u>	<u>20.7</u>	<u>19.7</u>	46.2	44.2	51.9	<u>47.9</u>	<u>38.1</u>	<u>41.1</u>	32.6	<u>59.9</u>	<u>25.3</u>	49.4	48.8	45.6	39.4
• SuMi [†] _{ICLR2025}	12.5	13.6	12.6	37.0	37.9	45.9	42.3	29.3	32.7	19.7	<u>55.7</u>	17.8	42.7	38.3	36.9	31.7
• BriMPR [†]	22.9	24.2	24.1	43.6	<u>45.4</u>	49.5	48.2	38.0	40.8	<u>36.8</u>	59.8	27.1	52.8	52.7	47.9	40.9

Table 3: Comparison with SOTA methods on Kinetics50-C under the multimodal shift setting (severity level 5).

Method	Noise			Weather			Avg.
	Gauss.	Traff.	Crowd	Rain	Thund.	Wind	
Source	17.1	6.4	5.4	6.0	13.5	8.8	9.5
• Tent	3.2	0.9	0.8	0.9	2.8	1.3	1.6
• EATA	21.5	7.7	7.1	7.3	17.3	11.9	12.1
• SAR	10.7	1.8	1.6	2.3	12.8	3.1	5.4
• DeYO	6.7	1.2	1.3	1.3	9.3	2.9	3.8
• FOA	18.8	10.8	11.4	<u>11.6</u>	20.5	10.4	13.9
• READ [†]	20.1	12.5	10.7	10.5	<u>20.5</u>	<u>13.4</u>	14.6
• ABPEM [†]	<u>21.9</u>	<u>13.4</u>	<u>12.3</u>	10.9	20.4	12.4	<u>15.2</u>
• SuMi [†]	17.0	6.8	5.7	6.2	13.4	8.8	9.7
• BriMPR [†]	23.5	18.8	21.4	15.8	26.8	18.3	20.7

Table 4: Comparison with SOTA methods on VGGSound-C under the multimodal shift setting (severity level 5).

domain shifts, we consider the settings of MOSI \rightarrow SIMS and SIMS \rightarrow MOSI, and report accuracy (ACC) and F1 score (F1).

Baselines. We compare the proposed method with multiple baselines including Source (source pre-trained model), Tent (Wang et al. 2021), EATA (Niu et al. 2022), SAR (Niu et al. 2023), DeYO (Lee et al. 2024), FOA (Niu et al. 2024), READ (Yang et al. 2024), ABPEM (Zhao et al. 2025) and

Method	MOSI \rightarrow SIMS		SIMS \rightarrow MOSI	
	ACC \uparrow	F1 \uparrow	ACC \uparrow	F1 \uparrow
Source	46.0	45.6	59.0	73.6
• Tent	38.1	42.2	59.6	74.5
• READ [†]	32.4	44.5	59.7	74.7
• SuMi [†]	44.4	45.0	59.4	74.2
• BriMPR [†]	58.2	57.6	59.9	74.9

Table 5: Comparison with SOTA methods on real-world shift datasets.

SuMi (Guo and Jin 2025).

Implementation details. For all experiments, we use an Adam optimizer with a learning rate of $1e-4$ and batch size of 64. The default number of prompts per layer m_p is set to 10 and the prompts are randomly initialized (Jia et al. 2022). The mask ratio is set to 0.5. τ_0 and D_0 of the adaptive temperature coefficient AdaTp are set to 0.2 and 5 respectively. τ in Eq. (9) is set to 0.07/0.25 for the unimodal and multimodal corruption settings respectively. For the hyperparameters of the compared methods, we adopt the recommended values from the respective papers. All the experiments are conducted with 3 random seeds on RTX-3090 GPUs.

Method	Kinetics50-C			VGGSound-C		
	audio	video	both	audio	video	both
BriMPR w/o $\mathcal{L}_{\text{CMER}}$	71.4	65.6	40.7	35.3	57.6	20.2
• BriMPR ($\lambda^a \leftrightarrow \lambda^v$)	70.0 (-1.4)	65.2 (-0.4)	39.9 (-0.8)	32.1 (-3.2)	56.5 (-1.1)	19.5 (-0.7)
• BriMPR	72.0 (+0.6)	65.9 (+0.3)	40.9 (+0.2)	36.5 (+1.2)	57.7 (+0.1)	20.7 (+0.5)

Table 6: Verify the effect of CMER from the perspective of weights.

Method	Kinetics50-C			VGGSound-C		
	audio	video	both	audio	video	both
Source	69.4	60.5	31.8	25.0	56.2	9.5
\mathcal{L}_{KL}	69.3	60.4	31.5	24.8	55.7	9.1
$\mathcal{L}_{\text{moment}_2}$	69.9	61.5	34.5	25.2	48.9	12.1
$\mathcal{L}_{\text{moment}_1}$	71.3	63.5	37.4	32.0	54.7	16.4
(A) $\mathcal{L}_{\text{PMGFA}}$	71.1	64.7	40.5	35.1	57.5	20.1
(B) + $\mathcal{L}_{\text{IICL}}$	71.4	65.6	40.7	35.3	57.6	20.2
(C) + $\mathcal{L}_{\text{CMER}}$	72.0	65.9	40.9	36.5	57.7	20.7

Table 7: Ablation studies for different components of BriMPR. \mathcal{L}_{KL} , $\mathcal{L}_{\text{moment}_2}$ and $\mathcal{L}_{\text{moment}_1}$ respectively denote replacing $\mathcal{L}_{\text{PMGFA}}$ with the KL-divergence, moment matching, and moment matching in a non-squared form.

Performance Comparison

Results of the unimodal shift setting. In Tab. 1 and Tab. 2, we present the results of the unimodal shift setting on Kinetics50-C and VGGSound-C with audio corruption and video corruption, respectively. Our proposed method BriMPR consistently improves the source model and outperforms all other competing methods. Notably, in scenarios where the dominant modality of the dataset is corrupted (for Kinetics50-C, video is the dominant modality; for VGGSound-C, audio is the dominant modality), BriMPR yields significant performance gains (60.5% \rightarrow 65.9% on Kinetics50-C; 25.0% \rightarrow 36.5% on VGGSound-C).

Results of the multimodal shift setting. Tab. 3 and Tab. 4 respectively present the results of the challenging multimodal shift setting on Kinetics50-C and VGGSound-C. Taking the ‘‘Gauss.’’ column in Tab. 3 as an example, the reported value denotes the average classification accuracy (%) across all 6 types of audio corruption, given the presence of Gaussian corruption in the video modality. Most methods suffer significant performance drops under this setting, whereas our BriMPR achieves the best results on most domains by decoupling MMTA into unimodal alignment subproblems, thereby reducing the dependence on high-quality modalities.

Results of the real-world shift setting. Tab. 5 presents the results from the MOSI/SIMS datasets using text, video, and audio modalities. BriMPR exhibits strong robustness to real-world shifts. Notably, only BriMPR achieves results better than random guess ($> 50\%$) on the MOSI \rightarrow SIMS task, thanks to its modulation of the target feature space.

Ablation Studies

Scrutinize CMER from the perspective of the weight λ^u .

To illustrate how multimodal test-time adaptation benefits from CMER, we swap the weights λ^u ($u \in \{a, v\}$) in $\mathcal{L}_{\text{CMER}}$, assigning lower weight to the augmentation with a milder distribution shift in the masked modality. As reported in Tab. 6, the mismatched weights lead to significant performance drops. Taking the case of audio corruption as an example (where $\lambda^v > \lambda^a$), the performance degradation can be attributed to two main factors: (1) For $\lambda^a \mathcal{L}_{avm}$, the small λ^a suppresses the extraction of multimodal information by the complete but low-quality audio modality; (2) For $\lambda^v \mathcal{L}_{amv}$, providing pseudo-labels to the augmentation with the masked audio modality introduces more error information into the unmasked high-quality video modality.

Component analysis. As shown in Tab. 7, we conducted an ablation study on the components of BriMPR. First, we verify the effectiveness of $\mathcal{L}_{\text{PMGFA}}$ (A); compared with KL-divergence (Row 2) and moment matching (Row 3), $\mathcal{L}_{\text{PMGFA}}$ demonstrates a significant advantage, as it eliminates the off-diagonal elements in the covariance matrix, reducing the estimation error. When moment matching is modified to a non-squared form (Row 4), performance improves in most cases, as the squared norm also amplifies the error. Subsequently, combining $\mathcal{L}_{\text{PMGFA}}$ (A), which serves as the initial alignment objective, with inter-modal instance-wise contrastive learning $\mathcal{L}_{\text{IICL}}$ (B) and cross-modal masked embedding recombination $\mathcal{L}_{\text{CMER}}$ (C) for alignment refinement, leads to further performance gains across all tasks.

Conclusion

In this paper, we introduce BriMPR, a novel MMTA method which tackles the coupling effect of unimodal feature shift and cross-modal semantic misalignment in a divide-and-conquer manner. Specifically, benefiting from the well-aligned source feature space, we first calibrate each unimodal global feature distribution via modality-specific prompts to achieve initial cross-modal semantic alignment. We then introduce a novel Cross-modal Masked Embedding Recombination strategy to facilitate the integration of multimodal information into low-quality modalities, and further refine the alignment via Inter-modal Instance-wise Contrastive Learning. Extensive experiments conducted on MMTA benchmark, which includes corruption datasets and real-world shift datasets, demonstrate the superiority of BriMPR over the SOTA methods.

Acknowledgments

This work was supported by the Fundamental Research Funds for the Central Universities (No. 2025JBZX059), the Natural Science Foundation of Hebei Province (No. F2025105018), the Tangshan Municipal Science and Technology Plan Project (No.23130225E) and the Beijing Natural Science Foundation (No.4242046).

References

- Chen, C.; Xie, W.; Huang, W.; Rong, Y.; Ding, X.; Huang, Y.; Xu, T.; and Huang, J. 2019. Progressive Feature Alignment for Unsupervised Domain Adaptation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Chen, H.; Xie, W.; Vedaldi, A.; and Zisserman, A. 2020a. Vggsound: A Large-Scale Audio-Visual Dataset. In *ICASSP 2020 - 2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 721–725.
- Chen, T.; Kornblith, S.; Norouzi, M.; and Hinton, G. 2020b. A Simple Framework for Contrastive Learning of Visual Representations. In *International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, 1597–1607. PMLR.
- Devlin, J.; Chang, M.-W.; Lee, K.; and Toutanova, K. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, 4171–4186. Minneapolis, Minnesota: Association for Computational Linguistics.
- Döbler, M.; Marsden, R. A.; and Yang, B. 2023. Robust Mean Teacher for Continual and Gradual Test-Time Adaptation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 7704–7714.
- Dong, H.; Chatzi, E.; and Fink, O. 2025. Towards Robust Multimodal Open-set Test-time Adaptation via Adaptive Entropy-aware Optimization. In *International Conference on Learning Representations*.
- Feng, C.-M.; Yu, K.; Liu, Y.; Khan, S.; and Zuo, W. 2023. Diverse Data Augmentation with Diffusions for Effective Test-time Prompt Tuning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2704–2714.
- Gan, Y.; Bai, Y.; Lou, Y.; Ma, X.; Zhang, R.; Shi, N.; and Luo, L. 2023. Decorate the Newcomers: Visual Domain Prompt for Continual Test Time Adaptation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, 7595–7603.
- Gong, Y.; Rouditchenko, A.; Liu, A. H.; Harwath, D.; Karlinsky, L.; Kuehne, H.; and Glass, J. R. 2023. Contrastive Audio-Visual Masked Autoencoder. In *International Conference on Learning Representations*.
- Guo, C.; Pleiss, G.; Sun, Y.; and Weinberger, K. Q. 2017. On Calibration of Modern Neural Networks. In *International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, 1321–1330. PMLR.
- Guo, Z.; and Jin, T. 2025. Smoothing the Shift: Towards Stable Test-Time Adaptation under Complex Multimodal Noises. In *International Conference on Learning Representations*.
- He, K.; Chen, X.; Xie, S.; Li, Y.; Dollár, P.; and Girshick, R. 2022. Masked Autoencoders Are Scalable Vision Learners. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 16000–16009.
- He, K.; Fan, H.; Wu, Y.; Xie, S.; and Girshick, R. 2020. Momentum Contrast for Unsupervised Visual Representation Learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Hinton, G.; Vinyals, O.; and Dean, J. 2015. Distilling the Knowledge in a Neural Network. arXiv:1503.02531.
- Jia, M.; Tang, L.; Chen, B.-C.; Cardie, C.; Belongie, S.; Hariharan, B.; and Lim, S.-N. 2022. Visual Prompt Tuning. In *European Conference on Computer Vision*, 709–727. Cham: Springer Nature Switzerland.
- Jin, Y.; Wang, X.; Long, M.; and Wang, J. 2020. Minimum Class Confusion for Versatile Domain Adaptation. In *European Conference on Computer Vision*, 464–480. Cham: Springer International Publishing.
- Kay, W.; Carreira, J.; Simonyan, K.; Zhang, B.; Hillier, C.; Vijayanarasimhan, S.; Viola, F.; Green, T.; Back, T.; Natsev, P.; Suleyman, M.; and Zisserman, A. 2017. The Kinetics Human Action Video Dataset. arXiv:1705.06950.
- Lee, J.; Jung, D.; Lee, S.; Park, J.; Shin, J.; Hwang, U.; and Yoon, S. 2024. Entropy is not Enough for Test-Time Adaptation: From the Perspective of Disentangled Factors. In *International Conference on Learning Representations*.
- Lester, B.; Al-Rfou, R.; and Constant, N. 2021. The Power of Scale for Parameter-Efficient Prompt Tuning. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, 3045–3059. Online and Punta Cana, Dominican Republic: Association for Computational Linguistics.
- Li, D.; Yang, Y.; Song, Y.-Z.; and Hospedales, T. 2018. Learning to Generalize: Meta-Learning for Domain Generalization. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32.
- Li, X. L.; and Liang, P. 2021. Prefix-Tuning: Optimizing Continuous Prompts for Generation. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, 4582–4597. Online: Association for Computational Linguistics.
- Li, Y.; Wang, N.; Shi, J.; Liu, J.; and Hou, X. 2017. Revisiting Batch Normalization For Practical Domain Adaptation. In *International Conference on Learning Representations*.
- Liu, X.; Ji, K.; Fu, Y.; Tam, W.; Du, Z.; Yang, Z.; and Tang, J. 2022. P-Tuning: Prompt Tuning Can Be Comparable to Fine-tuning Across Scales and Tasks. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, 61–68. Dublin, Ireland: Association for Computational Linguistics.

- Liu, Y.; Kothari, P.; van Delft, B.; Bellot-Gurlet, B.; Mordan, T.; and Alahi, A. 2021. TTT++: When Does Self-Supervised Test-Time Training Fail or Thrive? In *Advances in Neural Information Processing Systems*, volume 34, 21808–21820. Curran Associates, Inc.
- Long, M.; Cao, Y.; Wang, J.; and Jordan, M. 2015. Learning Transferable Features with Deep Adaptation Networks. In *International Conference on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*, 97–105. Lille, France: PMLR.
- Niu, S.; Miao, C.; Chen, G.; Wu, P.; and Zhao, P. 2024. Test-Time Model Adaptation with Only Forward Passes. In *International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, 38298–38315. PMLR.
- Niu, S.; Wu, J.; Zhang, Y.; Chen, Y.; Zheng, S.; Zhao, P.; and Tan, M. 2022. Efficient Test-Time Model Adaptation without Forgetting. In *International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, 16888–16905. PMLR.
- Niu, S.; Wu, J.; Zhang, Y.; Wen, Z.; Chen, Y.; Zhao, P.; and Tan, M. 2023. Towards Stable Test-time Adaptation in Dynamic Wild World. In *International Conference on Learning Representations*.
- Shin, I.; Tsai, Y.-H.; Zhuang, B.; Schuler, S.; Liu, B.; Garg, S.; Kweon, I. S.; and Yoon, K.-J. 2022. MM-TTA: Multi-Modal Test-Time Adaptation for 3D Semantic Segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 16928–16937.
- Shu, M.; Nie, W.; Huang, D.-A.; Yu, Z.; Goldstein, T.; Anandkumar, A.; and Xiao, C. 2022. Test-Time Prompt Tuning for Zero-Shot Generalization in Vision-Language Models. In *Advances in Neural Information Processing Systems*, volume 35, 14274–14289. Curran Associates, Inc.
- Su, Y.; Xu, X.; and Jia, K. 2022. Revisiting Realistic Test-Time Training: Sequential Inference and Adaptation by Anchored Clustering. In *Advances in Neural Information Processing Systems*, volume 35, 17543–17555. Curran Associates, Inc.
- Sun, Y.; Wang, X.; Liu, Z.; Miller, J.; Efros, A.; and Hardt, M. 2020. Test-Time Training with Self-Supervision for Generalization under Distribution Shifts. In *International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, 9229–9248. PMLR.
- Tzeng, E.; Hoffman, J.; Zhang, N.; Saenko, K.; and Darrell, T. 2014. Deep Domain Confusion: Maximizing for Domain Invariance. arXiv:1412.3474.
- Wang, D.; Shelhamer, E.; Liu, S.; Olshausen, B.; and Darrell, T. 2021. Tent: Fully Test-Time Adaptation by Entropy Minimization. In *International Conference on Learning Representations*.
- Wang, Q.; Fink, O.; Van Gool, L.; and Dai, D. 2022. Continual Test-Time Domain Adaptation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 7201–7211.
- Wang, Y.; Chauhan, J.; Wang, W.; and Hsieh, C.-J. 2023. Universality and Limitations of Prompt Tuning. In *Advances in Neural Information Processing Systems*, volume 36, 75623–75643. Curran Associates, Inc.
- Wang, Z.; Chi, Z.; Wu, Y.; Gu, L.; Liu, Z.; Plataniotis, K.; and Wang, Y. 2025. Distribution Alignment for Fully Test-Time Adaptation with Dynamic Online Data Streams. In *European Conference on Computer Vision*, 332–349. Cham: Springer Nature Switzerland.
- Yang, M.; Li, Y.; Zhang, C.; Hu, P.; and Peng, X. 2024. Test-time Adaptation against Multi-modal Reliability Bias. In *International Conference on Learning Representations*.
- Yoo, S.; Kim, E.; Jung, D.; Lee, J.; and Yoon, S. 2023. Improving Visual Prompt Tuning for Self-supervised Vision Transformers. In *International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, 40075–40092. PMLR.
- Yu, W.; Xu, H.; Meng, F.; Zhu, Y.; Ma, Y.; Wu, J.; Zou, J.; and Yang, K. 2020. CH-SIMS: A Chinese Multimodal Sentiment Analysis Dataset with Fine-grained Annotation of Modality. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 3718–3727. Online: Association for Computational Linguistics.
- Zadeh, A.; Zellers, R.; Pincus, E.; and Morency, L.-P. 2016. Multimodal Sentiment Intensity Analysis in Videos: Facial Gestures and Verbal Messages. *IEEE Intelligent Systems*, 31(6): 82–88.
- Zhang, J.; Huang, J.; Zhang, X.; Shao, L.; and Lu, S. 2024a. Historical Test-time Prompt Tuning for Vision Foundation Models. In *Advances in Neural Information Processing Systems*.
- Zhang, M.; Levine, S.; and Finn, C. 2022. MEMO: Test Time Robustness via Adaptation and Augmentation. In *Advances in Neural Information Processing Systems*, volume 35, 38629–38642. Curran Associates, Inc.
- Zhang, Y.; and Feng, S. 2023. Enhancing Domain-Invariant Parts for Generalized Zero-Shot Learning. In *Proceedings of the 31st ACM International Conference on Multimedia*, MM ’23, 6283–6291. New York, NY, USA: Association for Computing Machinery. ISBN 9798400701085.
- Zhang, Z.-Y.; Xie, Z.; Yao, H.; and Sugiyama, M. 2024b. Test-time Adaptation in Non-stationary Environments via Adaptive Representation Alignment. In *Advances in Neural Information Processing Systems*, volume 37, 94607–94632. Curran Associates, Inc.
- Zhao, Y.; Luo, J.; Luo, X.; Huang, J.; Yuan, J.; Xiao, Z.; and Zhang, M. 2025. Attention Bootstrapping for Multi-Modal Test-Time Adaptation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, 22849–22857.
- Zhou, K.; Yang, J.; Loy, C. C.; and Liu, Z. 2022a. Conditional Prompt Learning for Vision-Language Models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 16816–16825.
- Zhou, K.; Yang, J.; Loy, C. C.; and Liu, Z. 2022b. Learning to prompt for vision-language models. *International Journal of Computer Vision*, 130(9): 2337–2348.

Zhou, K.; Yang, Y.; Qiao, Y.; and Xiang, T. 2021. Domain Generalization with MixStyle. In *International Conference on Learning Representations*.

Appendix

Proof

Theorem 1. Given $x_1, \dots, x_n \in \mathbb{R}^d$ independently drawn from a multivariate normal distribution $\mathcal{N}(\mu, \Sigma)$, let $\hat{\Sigma}$ be the unbiased sample covariance matrix and $\hat{\sigma}^2 = [\hat{\sigma}_1^2, \dots, \hat{\sigma}_d^2]^T$ be the vector of its diagonal entries. Then, the mean squared errors satisfy:

$$\mathbb{E} [\|\hat{\Sigma} - \Sigma\|_F^2] = \mathcal{O}\left(\frac{d^2}{n}\right), \mathbb{E} [\|\hat{\sigma}^2 - \sigma^2\|_2^2] = \mathcal{O}\left(\frac{d}{n}\right). \quad (11)$$

Proof. The squared Frobenius norm of the error is:

$$\|\hat{\Sigma} - \Sigma\|_F^2 = \sum_{i=1}^d \sum_{j=1}^d (\hat{\Sigma}_{ij} - \Sigma_{ij})^2. \quad (12)$$

Taking expectations and using the unbiasedness of $\hat{\Sigma}$ ($\mathbb{E}[\hat{\Sigma}] = \Sigma$):

$$\begin{aligned} \mathbb{E} [\|\hat{\Sigma} - \Sigma\|_F^2] &= \sum_{i=1}^d \sum_{j=1}^d \mathbb{E} [(\hat{\Sigma}_{ij} - \Sigma_{ij})^2] \\ &= \sum_{i=1}^d \sum_{j=1}^d \text{Var}(\hat{\Sigma}_{ij}). \end{aligned} \quad (13)$$

Since $(n-1)\hat{\Sigma} \sim W_d(n-1, \Sigma)$, the Wishart distribution properties imply:

$$\text{Var}((n-1)\hat{\Sigma}_{ij}) = (n-1)(\Sigma_{ij}^2 + \Sigma_{ii}\Sigma_{jj}), \quad (14)$$

and thus:

$$\text{Var}(\hat{\Sigma}_{ij}) = \frac{\Sigma_{ij}^2 + \Sigma_{ii}\Sigma_{jj}}{n-1}. \quad (15)$$

Substituting Eq. (15) back into Eq. (13):

$$\begin{aligned} \mathbb{E} [\|\hat{\Sigma} - \Sigma\|_F^2] &= \frac{1}{n-1} \left(\sum_{i,j} \Sigma_{ij}^2 + \sum_{i,j} \Sigma_{ii}\Sigma_{jj} \right) \\ &= \frac{\|\Sigma\|_F^2 + (\text{tr}(\Sigma))^2}{n-1}. \end{aligned} \quad (16)$$

Assuming Σ has bounded entries (independent of n and d), $\|\Sigma\|_F^2 = \mathcal{O}(d^2)$ and $\text{tr}(\Sigma) = \mathcal{O}(d)$. Therefore:

$$\mathbb{E} [\|\hat{\Sigma} - \Sigma\|_F^2] = \mathcal{O}\left(\frac{d^2}{n}\right).$$

The squared L2 norm of the error is:

$$\|\hat{\sigma}^2 - \sigma^2\|_2^2 = \sum_{i=1}^d (\hat{\sigma}_i^2 - \sigma_i^2)^2. \quad (17)$$

Taking expectations and using $\mathbb{E}[\hat{\sigma}_i^2] = \sigma_i^2$:

$$\mathbb{E} [\|\hat{\sigma}^2 - \sigma^2\|_2^2] = \sum_{i=1}^d \mathbb{E} [(\hat{\sigma}_i^2 - \sigma_i^2)^2] = \sum_{i=1}^d \text{Var}(\hat{\sigma}_i^2). \quad (18)$$

Since $\hat{\sigma}_i^2 = \hat{\Sigma}_{ii}$, from the Wishart distribution:

$$\text{Var}(\hat{\sigma}_i^2) = \text{Var}(\hat{\Sigma}_{ii}) = \frac{\Sigma_{ii}^2 + \Sigma_{ii}\Sigma_{ii}}{n-1} = \frac{2\Sigma_{ii}^2}{n-1}. \quad (19)$$

Substituting Eq. (19) back into Eq. (18):

$$\mathbb{E} [\|\hat{\sigma}^2 - \sigma^2\|_2^2] = \sum_{i=1}^d \frac{2\Sigma_{ii}^2}{n-1} = \frac{2}{n-1} \sum_{i=1}^d \Sigma_{ii}^2. \quad (20)$$

Assuming Σ_{ii} are bounded, $\sum_{i=1}^d \Sigma_{ii}^2 = \mathcal{O}(d)$, yielding:

$$\mathbb{E} [\|\hat{\sigma}^2 - \sigma^2\|_2^2] = \mathcal{O}\left(\frac{d}{n}\right).$$

□

Algorithm of BriMPR

Algorithm 1: BriMPR

Input: test data stream $(\mathbf{x}_t^a, \mathbf{x}_t^v)$; modality-specific encoders Φ^a and Φ^v with prompts $\mathbf{P}^a = \{\mathbf{P}_0^a, \mathbf{P}_1^a, \dots, \mathbf{P}_{N-1}^a\}$ and $\mathbf{P}^v = \{\mathbf{P}_0^v, \mathbf{P}_1^v, \dots, \mathbf{P}_{N-1}^v\}$; joint module Ψ ; classification head h ; hyperparameters τ_0, D_0, τ .

- 1: **for** $t = 1$ to T **do**
 - 2: Get masked data $(\mathbf{x}^{a_m}, \mathbf{x}^{v_m}) \xleftarrow{\text{mask}} (\mathbf{x}^a, \mathbf{x}^v)$.
 - 3: Get fused features $\mathbf{Z}^{a_m v} = \Psi([\Phi^a(\mathbf{x}^{a_m}); \Phi^v(\mathbf{x}^{v_m})])$, $\mathbf{Z}^{a v_m} = \Psi([\Phi^a(\mathbf{x}^a); \Phi^v(\mathbf{x}^{v_m})])$ and $\mathbf{Z}^{a v} = \Psi([\Phi^a(\mathbf{x}^a); \Phi^v(\mathbf{x}^v)])$.
 - 4: Get augmented predictions $\mathbf{y}^{a_m v}, \mathbf{y}^{a v_m}$ via Eq. (5).
 - 5: Calculate $\text{Disc}^a, \text{Disc}^v$ and Disc^J .
 - 6: Calculate the adaptive temperature coefficient $\text{AdaTp} = 1 + \tau_0 / (1 + \exp(D_0 - \text{Disc}^J))$.
 - 7: Get calibrated pseudo-labels $\hat{\mathbf{y}}^{a v}$ via Eq. (6).
 - 8: Get unimodal features $\mathbf{Z}^a, \mathbf{Z}^v$ via Eq. (8).
 - 9: Calculate the overall loss $\mathcal{L}_{\text{BriMPR}}$ via Eq. (10).
 - 10: Update the modality-specific prompts \mathbf{P}^a and \mathbf{P}^v .
 - 11: **end for**
-

Experimental Details

Details of Baselines.

- **Tent** (Wang et al. 2021) optimizes the affine parameters in the normalization layer by minimizing the entropy of the model’s predictions. Because entropy can reflect the uncertainty of the predictions, and the normalization layer is associated with the distribution information. In the transformer-based CAV-MAE, we replace the Batch Normalization (BN) layers in the original implementation with Layer Normalization (LN) layers.
- **EATA** (Niu et al. 2022) selects reliable and on-redundant test samples with low entropy to participate in entropy minimization, and introduces a weighted Fisher regularizer to prevent significant changes in the parameters that are crucial for the in-distribution data during adaptation, thus alleviating the forgetting problem. The entropy threshold E_0 is set to $0.4 \times \ln C$ (where C is the number of classes). The cosine similarity threshold ϵ used for

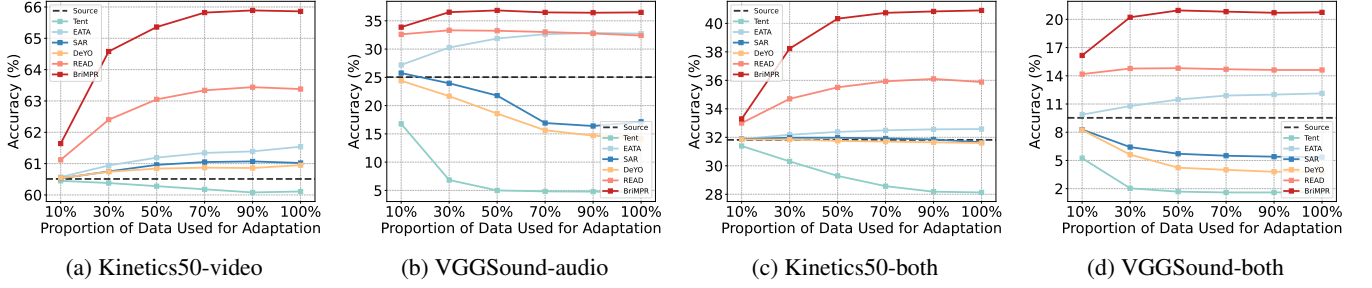


Figure 3: Performance comparison when the data available for adaptation is limited under two tasks on Kinetics50-C and VGGSound-C. (a) and (b) correspond to the unimodal shift setting; (c) and (d) correspond to the multimodal shift setting.

Method	Noise			Blur				Weather				Digital				Avg.
	Gauss.	Shot	Impul.	Defoc.	Glass	Motion	Zoom	Snow	Frost	Fog	Bright.	Contr.	Elast.	Pixel.	Jpeg	
Source	62.7	62.7	62.4	73.3	70.6	75.6	71.3	67.9	65.7	63.6	79.5	68.2	73.6	77.2	75.2	70.0
• Tent	63.1	63.0	62.9	73.4	71.3	<u>76.0</u>	72.1	68.6	66.6	63.7	79.3	68.4	74.3	77.1	75.0	70.3
• EATA	63.3	63.2	63.1	73.4	71.5	<u>76.0</u>	72.2	68.5	66.8	65.1	79.4	68.6	74.3	77.1	75.2	70.5
• SAR	63.2	63.0	62.8	73.2	71.5	76.1	72.4	68.5	67.0	65.1	79.4	68.7	74.2	77.1	75.1	70.5
• DeYO	63.3	63.3	63.2	73.4	71.4	75.9	72.2	68.5	66.7	64.5	79.3	68.7	74.2	77.1	75.3	70.5
• READ [†]	<u>64.0</u>	<u>63.9</u>	<u>64.1</u>	73.5	<u>72.1</u>	75.8	73.3	<u>69.5</u>	<u>68.8</u>	<u>68.9</u>	<u>79.5</u>	<u>69.8</u>	<u>74.9</u>	<u>77.3</u>	<u>75.7</u>	<u>71.4</u>
• SuMi [†]	63.1	63.2	63.1	73.6	71.5	75.9	72.1	68.7	66.7	64.4	79.6	68.7	74.2	<u>77.4</u>	75.5	70.5
• BriMPR [†]	67.2	67.1	67.0	73.6	73.7	75.8	<u>73.0</u>	71.0	70.0	70.9	79.2	71.1	75.6	77.8	76.1	72.6

Table 8: Comparison with SOTA methods on Kinetics50-C under the unimodal shift setting (mixed severity levels of video corruption).

Method	Noise			Weather			Avg.
	Gauss.	Traff.	Crowd	Rain	Thund.	Wind	
Source	76.7	65.0	68.5	69.2	68.9	72.1	70.1
• Tent	76.8	67.1	69.9	70.9	69.3	<u>72.8</u>	71.1
• EATA	<u>77.0</u>	67.4	70.1	<u>71.4</u>	70.6	72.5	71.5
• SAR	76.8	67.0	69.8	70.8	70.2	72.5	71.2
• DeYO	76.9	67.2	70.0	71.2	70.0	72.6	71.3
• READ [†]	77.1	<u>70.7</u>	<u>71.3</u>	72.5	<u>73.0</u>	72.5	<u>72.8</u>
• SuMi [†]	76.8	66.0	69.4	70.9	70.1	72.1	70.9
• BriMPR [†]	77.1	70.8	72.6	72.5	73.4	73.0	73.2

Table 9: Comparison with SOTA methods on Kinetics50-C under the unimodal shift setting (mixed severity levels of audio corruption).

filtering redundant samples is set to 0.1. The trade-off hyperparameter β is set to 1. The Fisher information is calculated using 2,000 unlabeled in-distribution samples, and the moving average factor α is set to 0.1.

- **SAR** (Niu et al. 2023) introduces sharpness-aware learning and minimizes entropy, optimizing the model weights to a flat minimum to enhance the robustness against noisy samples. Similar to EATA, the entropy threshold E_0 is set to $0.4 \times \ln C$. The radius ρ in the sharpness-aware optimization is set to 0.05. The model recovery threshold e_0 is set to 0.2. The moving average factor used to track the loss value is set to 0.9.
- **DeYO** (Lee et al. 2024) introduces the Pseudo-Label

Probability Difference (PLPD) metric to identify harmful samples that cannot be detected by entropy. It quantifies the impact of object shape information on predictions by measuring the difference in predictions before and after applying a single image transformation. The entropy threshold τ_{Ent} is set to $0.5 \times \ln C$, the PLPD threshold τ_{PLPD} is set to $0.2 \times \ln C$, and the normalization factor Ent_0 in the weighting function is set to $0.4 \times \ln C$. The default patch-shuffling is used as the image transformation.

- **FOA** (Niu et al. 2024) inserts prompts at the input level and employs the derivative-free covariance matrix adaptation (CMA) evolution strategy to learn the prompts. Additionally, it further adjusts the activation features at the output feature level to align them with the source domain. The number of prompt embeddings N_p is set to 1 with the default uniform initialization. The population size K in the CMA evolution strategy is set to $27 = 4 + 3 \ln(2 \times 768)$. The trade-off parameter λ in the fitness function is set to 0.4. The step size γ in the back-to-source activation shifting is set to 1.0. The moving average factor for computing the test statistics is set to 0.1.
- **READ** (Yang et al. 2024) updates the self-attention layer of the fusion module by optimizing the confidence-aware loss function, which promotes reliable fusion across different modalities. The confidence threshold γ is set to e^{-1} .
- **ABPEM** (Zhao et al. 2025) aligns cross-attention to self-

Method	Noise			Blur				Weather				Digital				Avg.
	Gauss.	Shot	Impul.	Defoc.	Glass	Motion	Zoom	Snow	Frost	Fog	Bright.	Contr.	Elast.	Pixel.	Jpeg	
Source	35.3	35.0	34.6	50.9	51.6	56.7	51.0	42.2	42.9	38.8	63.1	43.3	56.1	58.2	56.2	47.7
• Tent	30.4	29.9	29.8	47.1	50.8	55.9	50.1	38.9	41.2	32.5	63.3	39.1	56.7	56.8	54.7	45.2
• EATA	36.5	36.3	36.1	52.7	53.5	58.7	52.9	43.6	44.2	39.3	64.6	44.2	57.6	59.9	57.8	49.2
• SAR	35.2	34.9	34.4	52.1	53.0	58.1	52.9	43.4	44.1	37.1	63.8	42.4	57.1	59.1	56.9	48.3
• DeYO	34.5	33.9	33.7	52.1	53.1	58.4	52.7	43.0	44.2	37.4	64.3	42.8	57.5	59.2	56.8	48.2
• READ [†]	<u>39.4</u>	<u>39.3</u>	<u>38.9</u>	<u>57.3</u>	<u>56.0</u>	61.6	<u>56.0</u>	<u>47.2</u>	<u>48.5</u>	<u>46.2</u>	<u>66.5</u>	<u>47.8</u>	<u>59.6</u>	<u>62.6</u>	<u>60.3</u>	<u>52.5</u>
• SuMi [†]	35.3	35.1	34.8	51.5	52.5	57.7	52.0	42.6	43.5	37.9	63.9	43.2	56.9	58.8	56.7	48.2
• BriMPR [†]	45.2	44.5	44.4	<u>56.5</u>	57.8	<u>60.6</u>	57.2	49.4	50.3	51.1	66.8	51.8	61.2	64.4	62.4	54.9

Table 10: Comparison with SOTA methods on Kinetics50-C under the multimodal shift setting (mixed severity levels).

attention to reduce inter-modal differences, while excluding non-dominant class samples to reduce gradient noise in entropy loss. The threshold k for class ranking is set to 8/30 for Kinetics50-C and VGGSound-C. The weight λ for the attention bootstrapping loss is set to 1.

- **SuMi** (Guo and Jin 2025) progressively adapts to strongly out-of-distribution samples through interquartile range smoothing, selects high-quality samples via unimodal-assisted identification, and balances adaptation across modalities by leveraging mutual information sharing. The multimodal threshold γ_m and the normalization factor Ent_0 are set to $0.4 \times \ln C$. The unimodal threshold γ_u is set to e^{-1} . The smoothing coefficient β is set to 0.6/0.9, the weighting term λ is set to 5.0 and the unimodal assistance t is set to 1.0 by default for Kinetics50-C and VGGSound-C. For multimodal shift setting and real-world shift setting, we set the mutual information sharing term t_0 as $\text{iter}/2$.

Further Experiments

Limited data for adaptation. To explore the adaptation process of the proposed BriMPR, we further conduct experiments in the scenario where the data available for adaptation is limited. In this case, the model can only utilize a portion of the test data for adaptation in the initial stage and then make predictions on the remaining test data. This is because it is unrealistic for an intelligent perception system deployed in the real world to maintain updates all the time. An excellent multimodal test-time adaptation algorithm should still demonstrate good performance improvements even with limited target data and benefit from an increase in the available data. As shown in Fig. 3, BriMPR consistently maintains the best performance in all scenarios, demonstrating its data efficiency. Some methods even exhibit performance degradation with increased available test data (e.g., Tent, SAR, and DeYO in Fig. 3b), reflecting that the error accumulation (Chen et al. 2019) commonly faced by unimodal test-time adaptation is further exacerbated in multimodal scenarios.

Results under Mixed Severity Levels. Our main experimental results are obtained with the largest corruption severity level 5. We conduct experiments under both two settings with the mixed severity levels on Kinetics50-C to further

verify the robustness of our BriMPR. In this case, the severity level of the corrupted modality is randomly selected from 1 to 5. For each experiment, we generate three sets of test data with three random seeds. As shown in Tab. 8, Tab. 9, and Tab. 10, BriMPR significantly outperforms other methods in the vast majority of cases.

Results of Continual Multimodal Test-Time Adaptation. We further introduce continuously changing domains (Wang et al. 2022; Döbler, Marsden, and Yang 2023; Gan et al. 2023) and refer to this setting as Continual Multimodal Test-Time Adaptation (CMMTTA), which is more challenging because the model lacks domain labels during adaptation.

In the continual unimodal shift setting, only the corrupted modality undergoes distribution changes. In contrast, the continual multimodal shift setting involves distribution changes in one modality at a time. Taking “Kinetics50-both” as an example, each domain is represented as a combination of “video corruption + audio corruption”, where the modality changed is highlighted. Then the domain sequence is constructed as follows: Gauss. + Gauss. \rightarrow Gauss. + **Traff.** $\rightarrow \dots \rightarrow$ Gauss. + **Wind** \rightarrow **Shot.** + Wind \rightarrow Shot. + **Thund.** $\rightarrow \dots \rightarrow$ Shot. + **Gauss.** \rightarrow **Impul.** + Gauss. $\rightarrow \dots$, resulting in a total of $15 \times 6 = 90$ continuous domains.

To adapt BriMPR to the CMMTTA setting, we propose a simple variant, BriMPR-continual, which leverages the by-product unimodal distribution discrepancy Disc^u to detect domain shifts. Specifically, we maintain a sliding window of size w (e.g., $w = 10$) to store the most recent Disc^u values. For each new Disc_t^u , we calculate the Z-score with respect to the mean μ_t^u and standard deviation σ_t^u of the windowed values as: $Z_t^u = \frac{\text{Disc}_t^u - \mu_t^u}{\sigma_t^u}$. If $Z_t^u > k$ (e.g., $k = 5$), we interpret it as a significant domain shift for modality u , and accordingly re-initialize the modality-specific prompts \mathbf{P}^u . As shown in Fig. 4, all compared methods exhibit some degree of knowledge forgetting during continual adaptation. Among them, EATA (Niu et al. 2022) performs well after 20 domains due to its specially designed anti-forgetting mechanism. In contrast, our BriMPR-continual consistently improves over the Source baseline and achieves the best performance.

Efficiency comparisons. Tab. 11 compares the computational efficiency of different methods on VGGSound-C. Al-

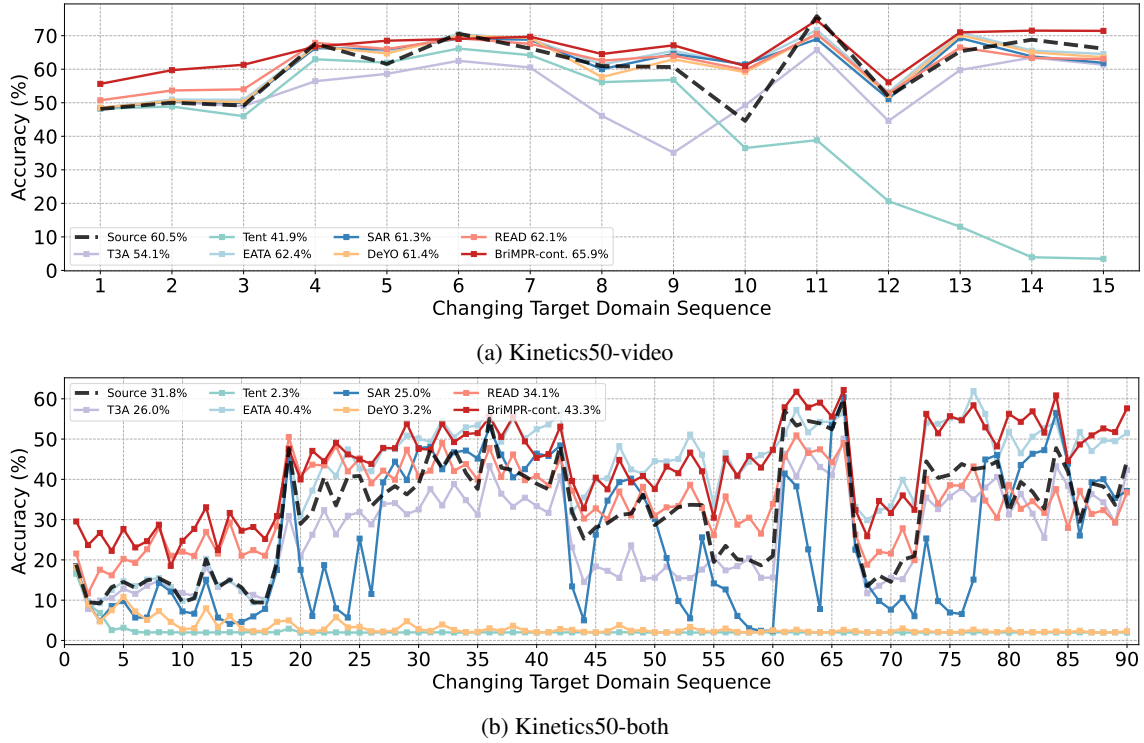


Figure 4: Comparison with the state-of-the-arts on Kinetics50-C under the continual setting (severity level 5). “Kinetics50-video” contains 15 continuous domains, while “Kinetics50-both” contains $15 \times 6 = 90$ continuous domains. The legend shows the average accuracy across all domains.

though BriMPR involves data augmentation and requires additional forward passes, the augmentation is performed via masking, making it more efficient than other augmentation-based methods like DeYO. Furthermore, thanks to the parameter efficiency of prompt tuning, BriMPR introduces fewer learnable parameters.

More Ablation Studies

Prompts are Better Distribution Calibrators. In many existing TTA works (Li et al. 2017; Wang et al. 2021; Niu et al. 2022), optimizing only the parameters of the normalization layer is regarded as a shortcut for calibrating the target feature distribution. In contrast, BriMPR calibrates the unimodal target feature distribution by optimizing the embedded prompts at each layer of the modality-specific encoders. As illustrated in Fig. 5, when using the same loss function $\mathcal{L}_{\text{PMGFA}}$, the variant $\mathcal{L}_{\text{PMGFA-LN}}$ that optimizes the LayerNorm parameters consistently underperforms prompt optimization across all tasks, while also requiring a larger number of trainable parameters.

Impact of τ_0 , D_0 and τ . We analyze the hyperparameters τ_0 , D_0 (in AdaTP) and τ (in Eq. (9)) on Kinetics50-C. As shown in Fig. 6a, τ_0 and D_0 are insensitive within reasonable ranges under both unimodal and multimodal corruption. This allows AdaTP to adaptively set temperatures and mitigate overconfident pseudo-labels. Fig. 6b shows that our method performs stably under different τ .

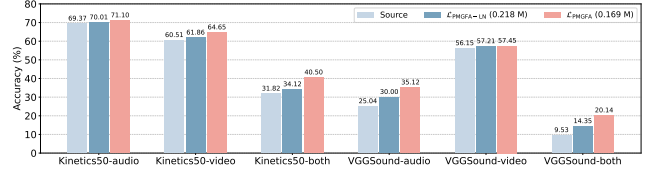


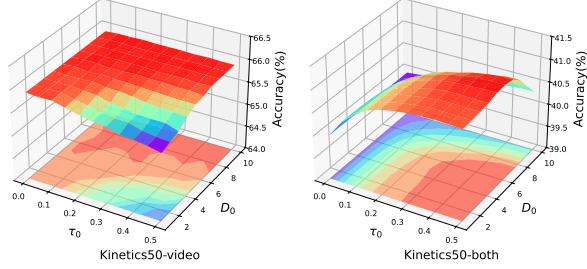
Figure 5: Comparison between updating LN parameters and updating prompts.

Impact of Number of Prompts. In Tab. 12, we present the results for different numbers of prompts, with 10 being the default value adopted by BriMPR. Previous prompt tuning work (Jia et al. 2022) indicates that the optimal number of prompts varies across different tasks. In our experiments on multimodal test-time adaptation task, adding more learnable prompts generally leads to performance improvements. Nevertheless, thanks to the parameter efficiency of prompt tuning, even a small number of prompts proves highly competitive, relieving us from the need to exhaustively search for the optimal prompt count.

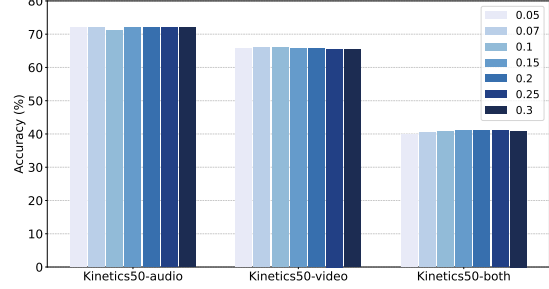
Impact of Number of Unlabeled Source Samples. By avoiding the use of covariance matrices in moment matching and adopting a non-squared formulation, $\mathcal{L}_{\text{PMGFA}}$ demonstrates greater robustness and significantly improves the quality of unimodal distribution alignment. In all

Method	Source	Tent	EATA	SAR	DeYO	READ [†]	SuMi [†]	BriMPR [†]
Time (s)	66.8	161.1	169.8	236.6	220.4	135.2	424.4	186.2
# Params (M)	0	0.218	0.218	0.218	0.218	1.772	0.218	0.169

Table 11: Efficiency comparisons among different methods on VGGSound-C.



(a) τ_0, D_0



(b) τ

Figure 6: Analysis of τ_0, D_0 and τ on Kinetics50-C.

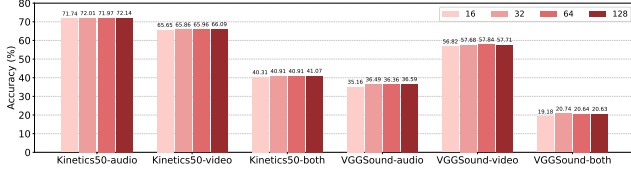


Figure 7: Comparison of different numbers of unlabeled source samples.

the experiments, we pre-estimate the source statistics $\{\hat{\mu}_i^{s,u}, \hat{\sigma}_i^{s,u}\}_{i=1}^N$ ($u \in \{a, v\}$) using 32 unlabeled source data. As shown in Fig. 7, BriMPR performs stably under different amounts of available source data.

Impact of Mask Ratio. In Fig. 8, we further investigate the sensitivity of BriMPR to the mask ratio in \mathcal{L}_{CMER} , traversing the values of [0.3, 0.4, 0.5, 0.6, 0.7] around the default value of 0.5. Our method maintains stable performance across different mask ratios. Meanwhile, it can be observed that: under the unimodal shift setting, increasing the mask ratio tends to improve performance (e.g., VGGSound-audio: 0.3 \rightarrow 0.7, accuracy 35.91% \rightarrow 36.97%); whereas under the multimodal shift setting, an excessively high mask ratio degrades performance (e.g., Kinetics50-both: 0.3 \rightarrow 0.7, accuracy 41.08% \rightarrow 40.19%). This is because under the former setting, there exists a clean modality, and boldly discarding the information of this modality promotes the recovery of the corrupted modality; whereas under the latter setting, the severity of shift varies among each modality, and excessive masking may suppress useful information and introduce noise from unreliable modalities.

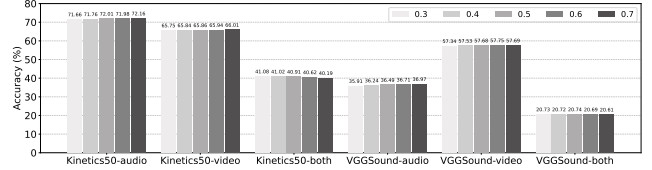


Figure 8: Comparison of different mask ratios.

Num. of Prompts	Ks50-audio	Ks50-video	Ks50-both	VGG-audio	VGG-video	VGG-both	Params (M)
1	71.22	63.80	37.95	35.58	57.62	19.05	0.017
3	71.48	64.55	39.24	36.31	57.74	20.30	0.051
5	71.73	64.91	39.94	36.33	57.69	20.48	0.084
10	72.01	65.86	40.91	36.49	57.68	20.74	0.169
20	72.19	66.07	41.90	36.50	57.69	20.98	0.338

Table 12: Comparison of different numbers of prompts.