# TRAINING-FREE TEST-TIME ADAPTATION WITH BROWNIAN DISTANCE COVARIANCE IN VISION-LANGUAGE MODELS

*Yi Zhang[1], Chun-Wun Cheng[2], Angelica I. Aviles-Rivero[3], Zhihai He[4], Liang-Jie Zhang[1]\**

[1]College of Computer Science and Software Engineering, Shenzhen University, China   [2]University of Cambridge, Cambridge, UK
[3]Yau Mathematical Sciences Center, Tsinghua University, Beijing, China   [4]Southern University of Science and Technology, Shenzhen, China

## ABSTRACT

Vision-language models suffer performance degradation under domain shift, limiting real-world applicability. Existing test-time adaptation methods are computationally intensive, rely on back-propagation, and often focus on single modalities. To address these issues, we propose *Training-free Test-Time Adaptation with Brownian Distance Covariance* (TaTa). TaTa leverages Brownian Distance Covariance—a powerful statistical measure that captures both linear and nonlinear dependencies via pairwise distances—to dynamically adapt VLMs to new domains without training or back-propagation. This not only improves efficiency but also enhances stability by avoiding disruptive weight updates. TaTa further integrates attribute-enhanced prompting to improve vision-language inference with descriptive visual cues. Combined with dynamic clustering and pseudo-label refinement, it effectively recalibrates the model for novel visual contexts. Experiments across diverse datasets show that TaTa significantly reduces computational cost while achieving state-of-the-art performance in domain and cross-dataset generalization.

***Index Terms***— Vision-Language Models, Test-Time Adaptation, Brownian Distance Covariance

## 1. INTRODUCTION

Vision-language models (VLMs) such as Align [1] and CLIP [2] learn aligned multimodal representations from large-scale image-text data, enabling effective zero-shot inference by comparing image and text embeddings in a shared space [3, 4, 5]. However, their performance degrades under significant domain or distribution shifts. Test-time adaptation (TTA) addresses this by adapting models using only unlabeled test data, without accessing source data or modifying training [6, 7]. This approach aligns well with practical scenarios where only the trained model is available, without access to the source data or authorization to alter the original training procedure. In such cases, models need to quickly adapt to new tasks.

Recent prompt-based TTA methods such as TPT [6] and DiffTPT [8] improve VLM adaptation but require computa-
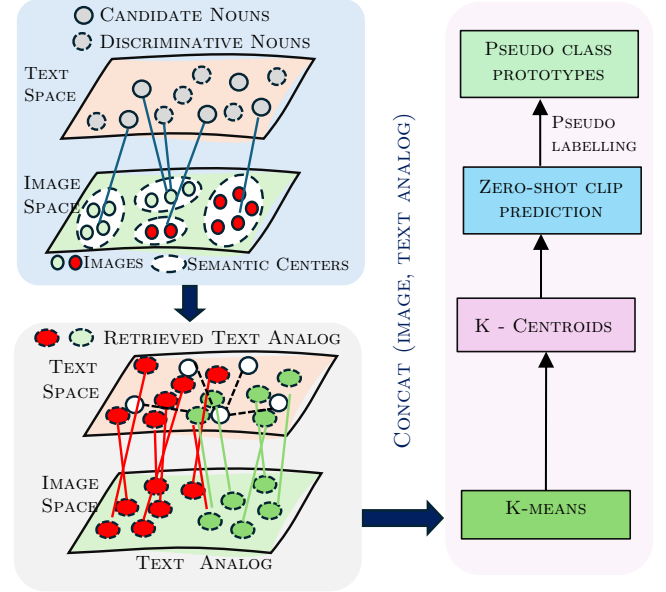


**Fig. 1**. The process involves using Pseudo-Labeling with Multimodal-Assisted Clustering to categorize nouns into semantic centers and create a discriminative text space. Nouns for each image are matched in this space, and by concatenating image and text features, K-means clustering forms centroids and clusters.

tionally expensive optimization [9]. TDA [10] avoids back-propagation and offers a training-free alternative, but has two key limitations: 1) TDA uses cosine similarity to measure the similarity between test image features and those in caches. This method only considers marginal distributions and linear relationships, potentially missing complex, non-linear dependencies in high-dimensional features. 2) TDA neglects the importance of visual attributes in prompts, relying solely on class labels or simple descriptions. This oversight limits the ability to capture the full semantic richness of images, resulting in less accurate and nuanced classifications, especially in complex or ambiguous scenarios.

To address these issues, we propose Training-free TTA method with Brownian Distance Covariance for VLMs. Brownian distance covariance can capture more comprehensive statistics by considering joint distribution. It effectively

---
*\*Corresponding author.*

models both linear and non-linear relationships between query and support images by measuring the discrepancy between the joint distribution of their features and the product of their marginals. The final prediction in our proposed TaTa is a combination of vision-vision and vision-language inference. For vision–vision inference, we use K-means with multimodal clues to cluster test images and assign pseudo-labels, as shown in Figure 1. BDC measures similarity between test features and pseudo-labeled centroids, capturing both linear and non-linear dependencies. For vision–language inference, we enrich prompts with visual attributes to improve semantic alignment. Finally, to address the prediction bias, we propose a soft-voting strategy that incorporates knowledge from nearest neighbor samples.

In summary, our contribution can be summarized as: 1) We present a streamlined and robust method for adapting vision-language models during testing without requiring additional training. 2) We introduce pseudo labeling with a dynamic clustering with multimodality, and adopt brownian distance covariance as the metrics for vision-vision inference, and we exploit visual attributes description words to enhance the prompting for vision-language inference. To mitigate pseudo-labels bias, we introduce a soft-voting strategy. 3) Extensive experimental results show that TaTa significantly surpasses existing state-of-the-art methods in both domain generalization and cross-dataset generalization tasks.

## 2. METHOD

**Method Overview.** Figure 2 illustrates our proposed TaTa framework for test-time adaptation, which enhances prediction by leveraging knowledge from the test stream and adapting image features. The final inference combines Vision–Vision(V-V) and Vision–Language(V-L) pathways. For a test image, we extract visual features using $E_v$ and combine them with textual analogs to form a multimodal feature $f_m$. In Vision–Vision Inference, we compute Brownian Distance Covariance (BDC) between $f_m$ and pseudo-labeled class prototypes derived from clustering. Simultaneously, for Vision–Language Inference, we construct attribute-enhanced prompts and compute cosine similarity between image features and corresponding text embeddings. A dynamic dictionary $\mathcal{D}$ is maintained, mapping class names to clusters. Correctly classified test samples update $\mathcal{D}$ with their multimodal features, continuously refining class centroids and reducing pseudo-label bias. A soft-voting mechanism incorporating nearest-neighbor knowledge further improves prediction robustness.

### 2.1. Dynamic Pseudo-Labeling with Multimodal-Assisted Clustering

A dynamic dictionary $\mathcal{D}$ stores class labels as keys and corresponding clusters as values. When a test image is correctly classified, its multimodal feature (concatenated image and text representations) is incorporated into the corresponding class in $\mathcal{D}$, updating the clusters and centroids.

**Dynamic Multimodal-Assisted Clustering:** To construct a discriminative textual space without known test class names, we use a subset of WordNet nouns [11]. As shown in Figure 1. For a dataset with $N$ classes, visual features extracted by $E_v$ are first clustered via K-Means to obtain initial centroids $C = \{C_i\}_{i=1}^N$. Using CLIP, we assign each WordNet noun to the closest semantic center based on similarity: $p(y = i \mid \mathbf{T_k}) = \frac{\exp(\text{sim}(f_{t_k}, C_i))}{\sum_{j=1}^N \exp(\text{sim}(f_{t_k}, C_j))}$, where $\mathbf{T_k}$ is a text prompt such as "a photo of noun" and $f_{t_k}$ is encoded by $E_t$. Top-$k_1$ nouns per center are selected to form a representative textual set $\{\bar{T}_m\}_{m=1}^H$ where $H = N \times k_1$ ($k_1 = 5$ empirically).

For each image feature $f_{v_h}$ of image $x_h$, we compute its textual analog $f_{t_h}$ by aggregating noun embeddings $\{\bar{f}_{t_h}\}_{h=1}^H$ with similarity-weighted coefficients: $f_{t_h} = \sum_{j=1}^H p(\bar{f}_{t_j} \mid f_{v_h})\bar{f}_{t_j}$, $p(\bar{f}_{t_j} \mid f_{v_h}) = \frac{\exp(\text{sim}(f_{v_h}, \bar{f}_{t_j})/\tilde{\tau})}{\sum_{l=1}^H \exp(\text{sim}(f_{v_h}, \bar{f}_{t_l})/\tilde{\tau})}$, where $\tilde{\tau} = 0.005$ (following TAC [12]). Finally, K-Means is applied to the concatenated features $[f_{t_d}, f_{v_d}]_{d=1}^D$ to obtain refined centroids $\bar{C} = \{\bar{C}_i\}_{i=1}^N$.

**Pseudo Labeling:** Text features are generated from class-based prompts using $E_t$. Pseudo-labels are assigned to each centroid $\bar{C}_i$ via zero-shot similarity comparison: $p(y = m \mid \bar{C}_i) = \frac{\exp(\text{sim}(f_{t_m}, \bar{C}_i)/\tau)}{\sum_{j=1}^N \exp(\text{sim}(f_{t_j}, \bar{C}_i)/\tau)}$, where $\tau$ is the softmax temperature. The resulting multimodal prototypes $\bar{C}$ serve as class representatives for inference.

### 2.2. BDC for V-V Inference

We design a Brownian Distance Covariance (BDC) module $\mathcal{B}(x)$ to capture both linear and nonlinear dependencies between feature representations. BDC measures dependence via distance covariance, based on pairwise distances between sampless [13, 14]. Given two random vectors $\mathbf{X}$ and $\mathbf{Y}$, the BDC is computed as follows: 1. Compute Euclidean distance matrices: $a_{ij} = \|\mathbf{X}_i - \mathbf{X}_j\|_2$, $b_{ij} = \|\mathbf{Y}_i - \mathbf{Y}_j\|_2$ 2. Center the distance matrices: $\bar{A}_{ij} = a_{ij} - \frac{1}{n}\sum_k a_{ik} - \frac{1}{n}\sum_l a_{lj} + \frac{1}{n^2}\sum_{k,l} a_{kl}$ (similarly for $\bar{B}_{ij}$) 3. Calculate the distance covariance: $\text{dCov}^2(\mathbf{X}, \mathbf{Y}) = \frac{1}{n^2}\sum_{i,j} \bar{A}_{ij}\bar{B}_{ij}$. BDC is non-zero only if features are dependent, capturing both linear and nonlinear relationships. It is parameter-free and training-free.

Inference: For each pseudo-labeled centroid $\bar{C}_i$ from clustering, we compute its BDC matrix $\mathcal{P}_{bdc}^i = \mathcal{B}(\bar{C}_i)$. For a test image $x_h$ with multimodal feature $f_m = [f_{v_h}, f_{v_t}]$, we compute its BDC matrix $f_{bdc}^m = \mathcal{B}(f_m)$. The prediction probability is:

$$p_{vv}(y = i|x_h) = \frac{\exp\left(\text{dCov}^2(f_{bdc}^m, \mathcal{P}_{bdc}^i)\right)}{\sum_{j=1}^N \exp\left(\text{dCov}^2(f_{bdc}^m, \mathcal{P}_{bdc}^j)\right)}.$$
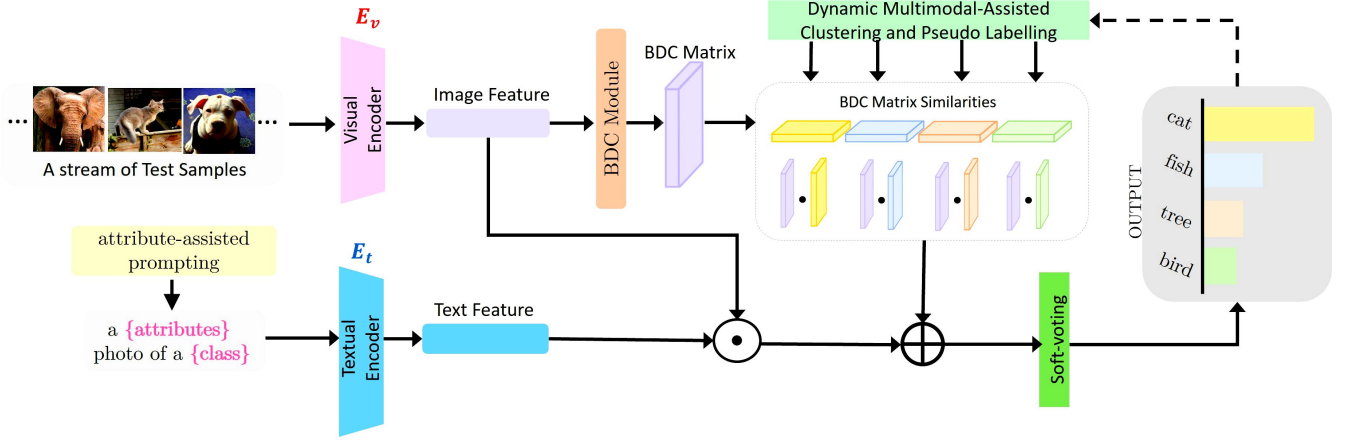
**Fig. 2**. An overview of our proposed TaTa. The final inference in TaTa combines Vision-Vision and Vision-Language Inference.

## 2.3. Attribute-assisted prompting for V-L Inference

Including descriptions of visual context can make the text more accurate, thereby improving the consistency between visual and linguistic modalities. This leads to better performance in CLIP's zero-shot inference. For example, "a photo of a grey koala on the tree" is more accurate than "a photo of a cat". Therefore, for the vision-language inference, we aim to exploit visual attributes to assist in prompting. First, we build a list $\Pi_t \triangleq \{\pi_i\}_{i=1}^{k_2}$ which contains $k_2 = 2000$ texts that describe the common visual attributes. The list $\Pi_t$ is derived from previous works [15, 16] and includes texts describing materials, patterns, colors, interactions, backgrounds, and more. Following CLIP's zero-shot setup, we concatenate each $\pi_i$ with a hand-crafted prompt $\psi$ "The photo is ..." to create a attribute-specific textual input $\{\psi; \pi_i\}$. This combined input is then fed into the text encoder $E_t$ to generate the attribute textual features $A_t \triangleq \{a_t^i\}_{i=1}^{k_2}$. For the image $x_{te}$, we use the image encoder $E_v$ to extract the image feature $f_v$, then we compute the cosine similarity with each attribute textual feature. We choose the attribute texts in $\Pi_t$ whose text features have the top-$k_2$ largest similarity score, to form the text "a {attributes} photo of a {class}". We then extract its feature $f_t$ using $E_t$ and perform the vision-language inference as:

$$p_{vl}(y=i|x_{te}) = \frac{\exp(\text{sim}(f_v, f_t^i/\tau)}{\sum_{j=1}^{N} \exp(\text{sim}(f_v, f_t^j)/\tau)}, \quad (1)$$

where $\tau$ is a temperature coefficient and "sim" denotes cosine similarity. $p_{vl}(y=i|x_{te})$ represents the prediction probability for $x_{te}$ belonging to label $i$.

## 2.4. Inference Fusion

We then combine $p_{vv}$ and $p_{vl}$ to get the final prediction,

$$p(y=i|x_{te}) = \alpha p_{vv}(y=i|x_{te}) + p_{vl}(y=i|x_{te}) \quad (2)$$

Prediction refinement is achieved by leveraging the knowledge from nearest neighbor samples, based on the premise that semantically similar images lie close in the feature space and are likely to share the same label. For a target image $x_{te}$, we extract its feature $f_v = E_v(x_{te})$, which is used to retrieve neighboring samples. The prediction for $x_{te}$ is then refined by aggregating the class probabilities of these neighbors via soft voting: $\hat{p}(y=i \mid x_{te}) = \frac{1}{k_3+1} \left( \sum_{j=1}^{k_3} p'_j + p(y=i \mid x_{te}) \right)$ where $k_3$ is the number of neighbors and $\hat{p}(y=i \mid x_{te})$ is the refined probability for class $i$.

## 3. EXPERIMENTS

### 3.1. Benchmark Settings

**Baselines.** In this paper, we comprehensively compare our proposed TaTa with seven state-of-the-art methods designed for vision-language models, including CLIP [2], CoOp [17], CoCoOp [18], Tip-Adapter [19], TPT [6], DiffTPT [8], and TDA [10]. Among these, CoOp, CoCoOp, and Tip-Adapter are train-time adaptation methods, while TPT, DiffTPT, and TDA are test-time adaptation methods. Train-time adaptation methods are trained on a 16-shot ImageNet set [20] and tested on other datasets, while test-time adaptation methods are fine-tuned directly on target datasets using the test set, without utilizing the ImageNet training set.

**Implementation details.** Our method utilizes ViT-B/16 CLIP for our experiments. Following prior methods [2], we adhere to data preprocessing protocols in CLIP. We empirically set $k_1 = 5$, $k_3 = 4$, and $\alpha = 1.75$. Following [10, 6], we report the top-1 accuracy (%), a standard classification criterion, as our evaluation metric.

### 3.2. Performance Analysis

**Domain Generalization.** In Table 2, we display the experimental results of the domain generalization benchmark, com-

**Table 1**. Comparison with state-of-the-art methods on the Cross-Domain Benchmark.

| Method | Aircraft | Caltech101 | Cars | DTD | EuroSAT | Flower102 | Food101 | Pets | SUN397 | UCF101 | Average |
|---|---|---|---|---|---|---|---|---|---|---|---|
| CLIP-ViT-B/16 | 23.22 | 93.55 | 66.11 | 45.04 | 50.42 | 66.99 | 82.86 | 86.92 | 65.63 | 65.16 | 64.59 |
| CoOp | 18.47 | 93.70 | 64.51 | 41.92 | 46.39 | 68.71 | 85.30 | 89.14 | 64.15 | 66.55 | 63.88 |
| CoCoOp | 22.29 | 93.79 | 64.90 | 45.45 | 39.23 | 70.85 | 83.97 | **90.46** | 66.89 | 68.44 | 64.63 |
| TPT | 24.78 | _94.16_ | 66.87 | _47.75_ | 42.44 | 68.98 | 84.67 | 87.79 | 65.50 | 68.04 | 65.10 |
| DiffTPT | _25.60_ | 92.49 | 67.01 | 47.00 | 43.13 | 70.10 | _87.23_ | 88.22 | 65.74 | 62.67 | 65.47 |
| TDA | 23.91 | **94.24** | _67.28_ | 47.40 | _58.00_ | **71.42** | 86.14 | 88.63 | _67.62_ | _70.66_ | _67.53_ |
| **TaTa (Ours)** | **26.68** | 93.82 | **68.04** | **49.57** | **59.67** | _71.28_ | **89.45** | _89.91_ | **68.64** | **73.54** | **69.06** |

**Table 2**. Comparison on domain generalization tasks.

| Method | ImageNet | -A | -V2 | -R | -S | OOD Average |
|---|---|---|---|---|---|---|
| CLIP | 68.34 | 49.89 | 61.88 | 77.65 | 48.24 | 59.42 |
| CoOp | **71.51** | 49.71 | 64.20 | 75.21 | 47.99 | 59.28 |
| CoCoOp | _71.02_ | 50.63 | 64.07 | 76.18 | 48.75 | 59.91 |
| Tip-Adapter | 70.75 | 51.04 | 63.41 | 77.76 | 48.88 | 60.27 |
| TPT | 68.98 | 54.77 | 63.45 | 77.06 | 47.94 | 60.81 |
| DiffTPT | 70.30 | 55.68 | 65.10 | 75.00 | 46.80 | 60.52 |
| TDA | 69.51 | _60.11_ | 64.67 | _80.24_ | _50.54_ | _63.89_ |
| **TaTa (Ours)** | 70.63 | **61.87** | **65.37** | **81.78** | **52.39** | **65.28** |

**Table 3**. Comparison of testing time and accuracy.

| Method | Testing Time | Accuracy | Gain |
|---|---|---|---|
| CLIP-ViT-B/16 | 12min | 68.34 | 0 |
| TPT | 12h 50min | 68.98 | +0.64 |
| TDA | 16min | 69.51 | +1.17 |
| **TaTa (Ours)** | **13.5min** | **70.63** | **+2.29** |

**Table 4**. Effectiveness of different components in TaTa. AAP refers to Attribute-Assisted Prompting. BDC represents the BDC module. MAC stands for Multimodal Assisted Clustering. SV denotes Soft-voting. The last row is our TaTa.

| Method | ImageNet | OOD Average |
|---|---|---|
| CLIP | 68.34 | 59.42 |
| + AAP | 69.05 | 61.21 |
| + AAP + BDC | 70.03 | 64.15 |
| + AAP + BDC +MAC | 70.42 | 64.85 |
| **+ AAP + BDC + MAC + SV** | **70.63** | **65.28** |

language models like CLIP, as it allows them to classify a wide range of classes in image classification without requiring additional training.

### 3.3. Ablation study

**Contributions of major algorithm components.** We conduct experiments on ViT-B/16 CLIP. As shown in Table 4, all three components significantly contribute to the overall performance improvement. Among them, the BDC module provides the largest performance gain, yielding a 2.94% improvement in the OOD average. This demonstrates that BDC provides an efficient metric for classification.

**Efficiency Comparison.** The comprehensive results are reported in Table 3. As shown, our method achieves new state-of-the-art performance in less time, exhibiting remarkable efficiency in test-time adaptation for VLMs. Compared to CLIP, TaTa gains 2.29% more accuracy with just 1.5 additional minutes. Compared to TDA, TaTa outperforms it by 1.12% and reduces testing time by 1.5 minutes.

### 4. CONCLUSION

We introduce TaTa, a training-free test-time adaptation method that addresses domain shift in vision-language models. By utilizing Brownian Distance Covariance and avoiding back-propagation, it reduces computational cost while maintaining stability. Through dynamic multimodal clustering and pseudo-labeling, TaTa enhances generalization performance across tasks. Experimental results confirm its effectiveness and practical potential, establishing TaTa as a robust solution for deploying adaptive vision-language models.

paring our proposed TaTa to various baselines. Compared to CLIP, TaTa shows significant improvements, with performance gains of up to 5.86% for OOD average. Notably, since train-time methods (CoOp and CoCoOp) are trained on the labeled training set of ImageNet, it is expected that they perform better on ImageNet. Interestingly, our method even outperforms Tip-Adapter on ImageNet, despite it being training-free but having access to the labeled training set. For the OOD average, TaTa consistently and substantially outperforms the train-time methods. When compared to test-time adaptation methods, our method outperforms all others by a large margin for both dataset-specific and OOD average results, except for ImageNet-V2. TaTa surpasses TDA by +1.39% on OOD average. These results demonstrate the effectiveness of our method and its strong test-time adaptation capability.

**Cross-dataset Generalization** Table 1 shows the comparison with state-of-the-art methods on the cross-dataset generalization task. It is obvious that our proposed TaTa substantially outperforms other methods on average. Compared to the second-best method TDA, our method achieves a performance gain of 1.53% . In particular, TaTa outperforms TDA by up to 2.77% on Aircraft and 2.88% on UCF101. Compared to train-time adaptation methods, TaTa surpasses them by up to 5.18% on average. These results highlight the strong test-time adaptation capability of our method across diverse class datasets. This feature is particularly beneficial for vision-

## 5. REFERENCES

[1] Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc Le, Yun-Hsuan Sung, Zhen Li, and Tom Duerig, "Scaling up visual and vision-language representation learning with noisy text supervision," in *International Conference on Machine Learning*. PMLR, 2021, pp. 4904–4916.

[2] 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, "Learning transferable visual models from natural language supervision," in *ICML*, 2021.

[3] Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee, "Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks," in *Advances in Neural Information Processing Systems*, 2019, vol. 32, pp. 13–23.

[4] Wonjae Kim, Bokyung Son, and Ildoo Kim, "Vilt: Vision-and-language transformer without convolution or region supervision," in *International conference on machine learning*. PMLR, 2021, pp. 5583–5594.

[5] Karan Desai and Justin Johnson, "Virtex: Learning visual representations from textual annotations," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 11162–11173.

[6] Manli Shu, Weili Nie, De-An Huang, Zhiding Yu, Tom Goldstein, Anima Anandkumar, and Chaowei Xiao, "Test-time prompt tuning for zero-shot generalization in vision-language models," *Advances in Neural Information Processing Systems*, vol. 35, pp. 14274–14289, 2022.

[7] Dian Chen, Dequan Wang, Trevor Darrell, and Sayna Ebrahimi, "Contrastive test-time adaptation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 295–305.

[8] Chun-Mei Feng, Kai Yu, Yong Liu, Salman Khan, and Wangmeng Zuo, "Diverse data augmentation with diffusions for effective test-time prompt tuning," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 2704–2714.

[9] Hee Suk Yoon, Eunseop Yoon, Joshua Tian Jin Tee, Mark A Hasegawa-Johnson, Yingzhen Li, and Chang D Yoo, "C-tpt: Calibrated test-time prompt tuning for vision-language models via text feature dispersion," in *The Twelfth International Conference on Learning Representations*, 2024.

[10] Adilbek Karmanov, Dayan Guan, Shijian Lu, Abdulmotaleb El Saddik, and Eric Xing, "Efficient test-time adaptation of vision-language models," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 14162–14171.

[11] George A Miller, "Wordnet: a lexical database for english," *Communications of the ACM*, vol. 38, no. 11, pp. 39–41, 1995.

[12] Yunfan Li, Peng Hu, Dezhong Peng, Jiancheng Lv, Jianping Fan, and Xi Peng, "Image clustering with external guidance," *arXiv preprint arXiv:2310.11989*, 2023.

[13] Gábor J Székely and Maria L Rizzo, "Brownian distance covariance," 2009.

[14] Jiangtao Xie, Fei Long, Jiaming Lv, Qilong Wang, and Peihua Li, "Joint distribution matters: Deep brownian distance covariance for few-shot classification," in *CVPR*, 2022.

[15] Yi Zhang, Ce Zhang, Ke Yu, Yushun Tang, and Zhihai He, "Concept-guided prompt learning for generalization in vision-language models," in *Proceedings of the AAAI Conference on Artificial Intelligence*, 2024, vol. 38, pp. 7377–7386.

[16] Yi Zhang, Ke Yu, Siqi Wu, and Zhihai He, "Conceptual codebook learning for vision-language models," *arXiv preprint arXiv:2407.02350*, 2024.

[17] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu, "Learning to prompt for vision-language models," *International Journal of Computer Vision*, vol. 130, no. 9, pp. 2337–2348, 2022.

[18] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu, "Conditional prompt learning for vision-language models," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 16816–16825.

[19] Renrui Zhang, Wei Zhang, Rongyao Fang, Peng Gao, Kunchang Li, Jifeng Dai, Yu Qiao, and Hongsheng Li, "Tip-adapter: Training-free adaption of clip for few-shot classification," in *European Conference on Computer Vision*. Springer, 2022, pp. 493–510.

[20] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "ImageNet: A Large-Scale Hierarchical Image Database," in *CVPR09*, 2009.