

# All Centers Are at most a Few Tokens Apart: Knowledge Distillation with Domain Invariant Prompt Tuning

Amir Mohammad Ezzati, Alireza Malekhosseini, Armin Khosravi, and  
Mohammad Hossein Rohban

Department of Computer Engineering, Sharif University of Technology, Tehran, Iran  
iamirezzati@gmail.com, alrezmlk@gmail.com, arminkhosravie@gmail.com,  
rohban@sharif.edu

**Abstract.** Domain generalization is critical in computational pathology (CPath) due to inherent domain shifts caused by variations in staining protocols, scanner devices, and imaging settings across clinical centers. Vision-language models (VLMs), such as PLIP—a pathology-tuned CLIP—trained on image-text pairs across diverse domains, serve as strong knowledge distillation sources. However, their zero-shot performance with predefined prompts remains limited due to sensitivity to prompt variations. Moreover, unlike natural images, histopathology centers lack semantic descriptors (e.g., 'sketch'), making it difficult to define domain-specific prompts for clinical centers. This requires a data-driven approach for learning domain-specific and ultimately class-generic continuous prompts. We propose Domain Invariant Prompt Tuning (DIPT) for knowledge distillation process, a novel step that learns multiple input tokens for each domain. These tokens are trained separately for each domain and are averaged across domains, leading to domain-invariant prompts. Our student model then distills knowledge from PLIP’s text encoder by leveraging the prompts learned by DIPT. This leads to alignment of visual features with domain-invariant embeddings, enhancing generalization by training on multiple domains. Our method adds a significant improvement in average F1-score to existing state-of-the-art (SOTA) knowledge distillation approaches in domain generalization with histopathology datasets. This work helps the way of deploying robust CPath models in real-world clinical problems with heterogeneous data sources. The code is available at [github.com/amirezzati/dipt](https://github.com/amirezzati/dipt).

**Keywords:** Domain generalization · Prompt tuning · Knowledge distillation · Vision Language Model · Computational pathology.

## 1 Introduction

Deep learning has demonstrated exceptional effectiveness in Computational Pathology (CPath), enabling accurate histology image classification. However, domain shift reduces model generalization to unseen datasets, necessitating domain generalization (DG) solutions [7].

Domain shift is a critical challenge in digital pathology, arising from differences in slide preparation, staining protocols, and scanner properties across medical centers. Even within the same institution, changes in the imaging pipeline over time can affect domain statistics. This issue results in deep learning models performing well on the training distribution but significantly worse on unseen domains [14,10].

Vision-language models (VLMs) like CLIP [12] have recently shown great promise in enhancing out-of-distribution/domain generalization in the zero-shot setup [3]. They have demonstrated adaptability to various downstream tasks by learning joint embeddings of images and text through contrastive learning [12]. Various follow up methods tried to enhance these models further through few-shot adaptation, showing remarkable robustness to domain drift [13].

Despite all these advancements, CLIP is trained mostly on natural images, limiting its applicability in the domain of medical images. To address this challenge, Pathology language-image pretraining (PLIP), a CLIP-based model fine-tuned on OpenPath dataset, was introduced. It demonstrates strong zero-shot performance on unseen pathology domains corresponding to imaging centers [6].

While pathology-pretrained models such as PLIP demonstrate remarkable zero-shot generalization across diverse domains, their absolute accuracy particularly for novel, specialized CPath tasks, fails to meet the necessary thresholds required for clinical deployment. One recent approach for addressing this issue is by improving domain generalization through Knowledge Distillation (KD) [1]. Specifically, while naive KD struggles to enhance out-of-distribution (OOD) generalization, integrating distillation from other modalities, such as text, into the student model improves generalization.

Aligned with this idea, some recent methods [1,5] in natural images propose distilling VLMs’ text encoders into the vision student model to enhance domain generalization. This is based on assuming text representations vary only slightly across domains, and are more generalizable than image representations [5,1]. These methods usually append some descriptor of the domain into the text input, e.g. “art” for an artistic domain to account for slight deviations from the class-generic representation.

However, applying such methods in CPath is challenging due to the absence of semantic descriptors (e.g., “art” or “sketch”) for histopathological data and the fact that existing class-generic prompts are not well-defined for histopathological images. In CPath, domains usually correspond to various imaging centers, which cannot be described expressively similar to natural images. To address this, we introduce Domain Invariant Prompt Tuning (DIPT), a novel step that generates domain-invariant and domain-specific prompts, enabling more effective knowledge distillation from PLIP’s text encoder.

Specifically, our method consists of two main steps. In the first step, we concat learnable tokens alongside a fixed aggregated generic token to capture unknown domain-specific knowledge. These learnable tokens are trained using prefix tuning, allowing them to adapt to domain-specific variations while retaining generalizable features. After training, we compute embeddings of all domain-

specific prompts for each class by passing them through the PLIP text encoder. By aggregating these domain-specific embeddings for a given class, we derive a domain-invariant class-generic embedding, which serves as a more general representation across domains.

In the second step, we freeze these domain-invariant class-generic embeddings for all classes and apply a KD pipeline to adapt the student model to various domains. This is achieved by distilling knowledge from both class-generic embeddings and the PLIP image encoder, ensuring the student model learns a more robust and generalizable representation. Experimental results on various benchmarks demonstrate the advantage of our model over zero-shot PLIP and other baselines, highlighting its effectiveness in improving domain generalization for computational pathology. Our contributions can be summarized as follows:

- We are the first to apply knowledge distillation from vision-language models like PLIP to enhance domain generalization in histopathology datasets.
- We introduce a novel step, DIPT, before knowledge distillation to generate domain-specific and well-defined class-generic prompts.
- We conduct extensive evaluations across various source and target domain combinations, comparing different knowledge distillation (KD) approaches on the **CAMELYON17-WILDS** [9] and **Kather19** [8] datasets, achieving improvements in the F1-score metric of up to **7.7%** compared to baseline methods.

## 2 Backgrounds

**Distillation from Vision Language Models.** Knowledge distillation (KD) [4] transfers knowledge from a high-capacity teacher model to a smaller student model, enabling the student to mimic the teacher’s outputs efficiently. Early KD methods relied on soft labels from teacher logits, which improve performance over direct supervision [4,11,15]. Later, H. Chen et al. introduced feature embedding distillation [2], which leverages embeddings rather than logits.

For vision-language models (VLMs), KD benefits from an additional source: the text encoder. Traditional KD focused on vision models, transferring visual representations to a student model. However, recent research shows that text encoders provide richer, domain-invariant semantic knowledge compared to image embeddings.

This shift has led to a new KD paradigm in VLMs, where knowledge transfer extends beyond visual representations. Notably, RISE [5] and VL2V-ADiP [1] leverage text encoders to enhance vision student models with more generalizable information.

Both methods employ a cosine similarity loss (see Eq. 1) to align student image embeddings with those of the teacher model, similar to traditional distillation methods. This encourages the student’s feature space to mimic the generalizable feature space of CLIP’s image encoder, which has been trained on large-scale datasets using contrastive learning.

$$\mathcal{L}_I = \sum_{(\mathbf{x}, \mathbf{y})} \cos(f(\mathbf{x}), h_I(\mathbf{x})), \quad (1)$$

where  $f(\cdot)$  and  $h_I(\cdot)$  are the student and teacher image encoders, respectively, and  $\cos(\cdot, \cdot)$  refers to the cosine similarity function.

RISE introduces a novel loss function that distills knowledge from CLIP’s text encoder. This loss (see Eq. 2) measures the absolute distance between the student’s image embedding and a generic text representation from the teacher’s text encoder, which is obtained by aggregating predefined template prompts across all domains. This alignment enhances knowledge transfer from CLIP, improving generalization.

$$\mathcal{L}_A = \sum_{(\mathbf{x}, \mathbf{y})} \sum_{i \in C} \mathbf{1}_{[y=i]} \cdot \cos(f(\mathbf{x}), \mathbf{E}_i). \quad (2)$$

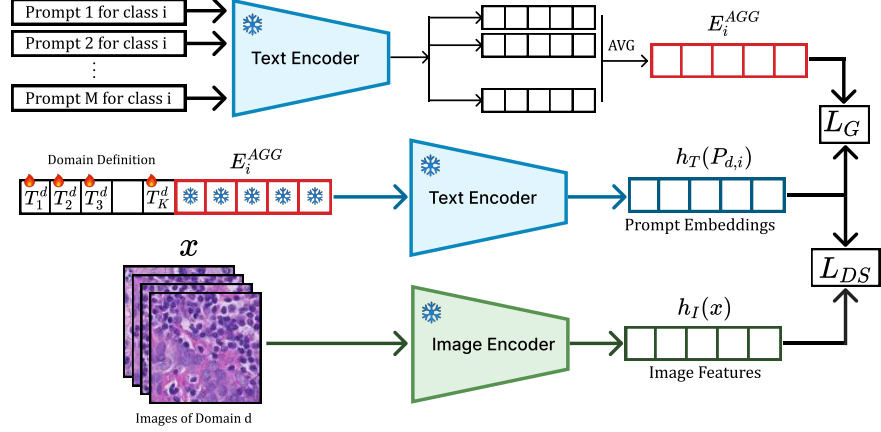
Here,  $\mathbf{E}_i$  is a domain-invariant class-generic embedding for class  $i$ ,  $\mathbf{1}$  is an indicator function, and  $C$  denotes the set of all possible classes.

However, in CPath, there is no well-defined, class-generic prompt that can be directly used as  $\mathbf{E}_i$  in the learning process. Moreover, there is no domain-specific descriptor for each domain, as the domains are different medical imaging centers, making it difficult to construct a domain-specific prompt that could be aggregated into a class-generic prompt.

**Prompt Tuning.** Prompt tuning adapts pretrained VLMs like CLIP by introducing task-specific textual tokens. While hand-crafted prompts (e.g., “a photo of a [CLASS]”) enable zero-shot classification, they often lack the flexibility to adapt to task-specific details. Context Optimization (CoOp) [19] replaces hand-crafted prompts with learnable soft prompts, improving adaptability. To further enhance generalization, some approaches [18, 17] introduce image-conditioned prompts, where image features are integrated with trainable textual prompts. Additionally, Knowledge-Guided Context Optimization (KgCoOp) [16] refines these techniques by addressing the problem of overfitting to seen classes. While CoOp boosts accuracy on seen classes, it often struggles with unseen classes, partly due to the gap between learned prompts and CLIP’s generic prompt embeddings. KgCoOp mitigates this issue by incorporating a knowledge-guided loss function, which ensures learned prompts retain general linguistic knowledge while being optimized for specific tasks.

### 3 Proposed Method

While [16] focuses on generalization to unseen classes, our goal is to improve generalization across unseen domains in CPath. To achieve this, we introduce DIPT (Domain-Invariant Prompt Tuning), a novel preparatory step before knowledge distillation. DIPT extends the pipeline of knowledge distillation by first generating *domain-specific* prompts through adaptation on the PLIP model. These



**Fig. 1. Domain-Specific Prompt Learning in DIPT.** Class-generic tokens are aggregated embeddings from multiple generic prompts per class (encoded via PLIP’s text encoder). These tokens are concatenated with K learnable continuous tokens used for adopting domain characteristics. Training these tokens adopts learnable prompts on PLIP for better classification.

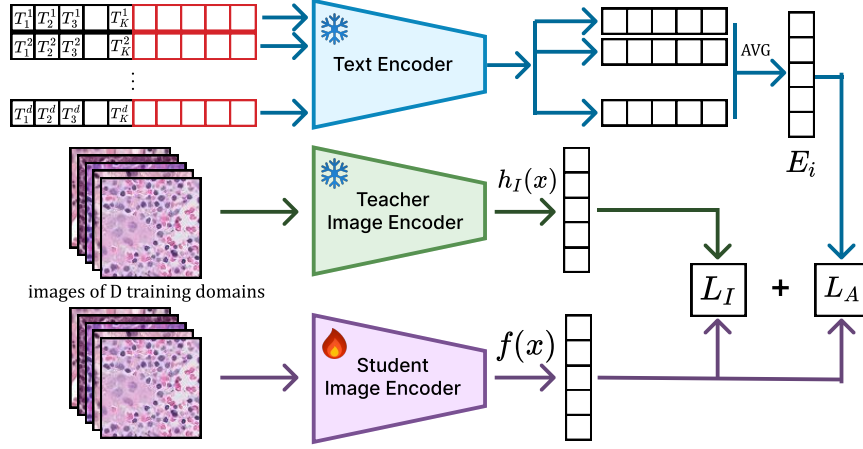
domain-specific prompts are then aggregated to form domain-invariant class-generic embeddings. Applying DIPT before PLIP knowledge distillation demonstrates stronger generalization and robustness compared to pure KD approaches.

### 3.1 Domain-Specific Prompt Learning

The first step of DIPT i.e domain-specific prompt learning optimizes learnable tokens to capture domain-relevant features while maintaining generalizable knowledge through aggregated template embeddings. There are two types of tokens: domain-specific tokens and class-generic tokens, with the final prompt formed by their concatenation (as shown in Fig. 1).

Domain-specific tokens are  $k$  learnable vectors that adapt to domain-specific semantics. These tokens act as histopathological domain descriptors, which do not exist beforehand but are learned to represent domain characteristics such as staining and imaging device variations through feedback from the image encoder ( $\mathcal{L}_{DS}$  loss in Fig. 1). These tokens are initialized using a Gaussian distribution and updated during training.

Class-generic tokens are frozen vectors that serve as a stable reference for each class and come from aggregated template embeddings. To construct these template embeddings, we generate  $M$  prompt templates for each class. For example, for the class “normal lymph node,” prompts like “a patch of normal lymph node” and “benign lymphoid cells” are used. These prompts are encoded using PLIP’s text encoder, and their embeddings are averaged per class to form the final aggregated template embeddings.



**Fig. 2. Knowledge Distillation via Class-Generic Learned Prompts.** Domain-specific prompts from  $D$  domains are aggregated after passing through PLIP’s text encoder  $h_T(\cdot)$  to form domain-invariant class-generic embeddings  $\mathbf{E}_i$  (Eq. 8). These embeddings are then used for dual knowledge distillation.

Formally, for the  $i$ -th class, the aggregated template embedding  $\mathbf{E}_i^{Agg}$  is computed as:

$$\mathbf{E}_i^{Agg} = \frac{1}{M} \sum_{m=1}^M h_T(\mathbf{P}_{i,m}^{Template}), \quad (3)$$

where  $h_T(\cdot)$  is the PLIP text encoder,  $\mathbf{P}_{i,m}^{Template}$  is the  $m$ -th prompt template for class  $i$ . During training, domain-specific learnable tokens  $T_1^d, T_2^d, \dots, T_K^d$  are optimized using a combined loss function that balances domain-specific discrimination and alignment with class-generic knowledge. The total loss (see Eq. 7) integrates a standard cross-entropy term as domain-specific loss:

$$\mathcal{L}_{DS} = - \sum_{(\mathbf{x}, \mathbf{y})} \log \frac{\exp(z_y/\tau)}{\sum_{j=1}^{N_c} \exp(z_j/\tau)}, \quad (4)$$

where,

$$z_i = \cos(h_I(\mathbf{x}), \mathbf{E}_{d,i}), \quad (5)$$

and a generalization loss term:

$$\mathcal{L}_G = \frac{1}{N_c} \sum_{i=1}^{N_c} \cos(\mathbf{E}_{d,i}, \mathbf{E}_i^{Agg}), \quad (6)$$

where  $\mathbf{E}_{d,i}$  denotes the domain-specific prompt embedding for domain  $d$  and class  $i$ .  $N_c$  is the total number of classes. Additionally,  $\tau$  is a temperature parameter,

and  $z_i$  represents the logit of class  $i$ . These two losses are added to form the final loss:

$$\mathcal{L} = \mathcal{L}_{DS} + \mathcal{L}_G. \quad (7)$$

### 3.2 Knowledge Distillation via Class-Generic Learned Prompts

To improve cross-domain adaptability, we aggregate  $D$  domain-specific prompts learned in the previous step into a domain-invariant class-generic embedding (see Eq. 8). This enhances generalization by capturing stronger cross-domain semantics while reducing domain-specific biases, leading to better performance on both seen and unseen domains:

$$\mathbf{E}_i = \frac{1}{D} \sum_{d=1}^D h_T(\mathbf{P}_{d,i}), \quad (8)$$

where  $h_T(\cdot)$  is the PLIP’s text encoder,  $\mathbf{P}_{d,i}$  is the domain-specific prompt for class  $i$  in the domain  $d$ .

Finally, this class-generic embedding  $\mathbf{E}_i$  is incorporated into the distillation pipeline (as shown in Fig. 2). This process helps the model develop a broader understanding that extends beyond individual domains. By leveraging well-defined class-generic embeddings, the model enhances its ability to generalize across diverse medical imaging centers, improving generalization to unseen domains.

**Table 1.** Test result on Camelyon17-WILDS. The results are reported with mean error bars, where the accuracy has a mean error of  $\pm 0.18\%$  and the F1 score has a mean error of  $\pm 0.21\%$ . The subscript \* indicates training with a ViT-based student, while others use ResNet-50. All reported values are in percent. (best results in **bold**)

Method	Center 1		Center 3		Center 4		Center 5	
	ACC	F1	ACC	F1	ACC	F1	ACC	F1
Agg. prompt Zero-Shot	81.89	81.79	89.38	89.64	77.63	77.27	75.58	79.98
KD [4]	93.94	93.61	88.57	87.39	91.42	90.75	83.32	80.20
RISE [5]	90.16	89.16	<b>93.07</b>	<b>92.85</b>	90.51	89.58	81.67	77.75
RISE + <b>DIPT</b>	<b>95.98</b>	<b>95.87</b>	91.23	90.51	<b>94.13</b>	<b>93.84</b>	<b>83.37</b>	<b>80.14</b>
VL2V [1]	93.24	92.76	90.43	89.47	<b>90.99</b>	<b>90.11</b>	87.42	85.78
VL2V + <b>DIPT</b>	<b>93.85</b>	<b>93.46</b>	<b>93.87</b>	<b>93.55</b>	90.90	89.99	<b>93.85</b>	<b>93.46</b>
VL2V* [1]	96.32	96.31	93.81	93.49	<b>95.24</b>	<b>95.06</b>	88.66	87.23
VL2V* + <b>DIPT</b>	<b>96.85</b>	<b>96.86</b>	<b>96.45</b>	<b>96.40</b>	94.47	94.18	<b>93.66</b>	<b>93.28</b>

## 4 Experiments and Results

**Datasets.** The Camelyon17-WILDS dataset [9] consists of tissue patches from five different hospitals, focusing on breast cancer metastases in lymph nodes.

Additionally, the Kather19 dataset [8] consists of colorectal tissue patches collected from three different centers, covering nine tissue types for histopathology analysis. They evaluate model generalization across unseen domains.

**Prompt Learning and Knowledge Distillation.** In both steps, we reserve one of the centers (e.g. center 2 for Camelyon17) for validation. We experiment with  $k \in \{2, 3, 4\}$  and various learning rates (e.g.  $5 \times 10^{-6}$ ,  $5 \times 10^{-5}$ ). For knowledge distillation, we evaluate VL2V (with ResNet-50 and ViT-B/16) and RISE (with ResNet-50) on Camelyon17 using a cross-domain setup that rotates the remaining four domains between training and test sets.

**Table 2.** Mean and worst-case performance from Table 1 on Camelyon17 (best results in **bold**).

Method	Mean		Worst	
	ACC	F1	ACC	F1
Zero-Shot	81.12	82.17	75.58	77.27
KD [4]	89.31	87.98	83.32	80.20
RISE [5]	88.85	87.33	81.67	77.75
RISE + <b>DIPT</b>	<b>91.17</b>	<b>90.09</b>	<b>83.37</b>	<b>80.14</b>
VL2V [1]	90.52	89.53	87.42	85.78
VL2V + <b>DIPT</b>	<b>93.11</b>	<b>92.61</b>	<b>90.90</b>	<b>89.99</b>
VL2V* [1]	93.50	93.02	88.66	87.23
VL2V* + <b>DIPT</b>	<b>95.35</b>	<b>95.18</b>	<b>93.66</b>	<b>93.28</b>

**Table 3.** Test result comparison on Kather19 (best results in **bold**).

Method	ACC	F1
Zero-Shot	63.66	60.86
KD [4]	92.87	92.98
RISE [5]	92.98	92.45
RISE + <b>DIPT</b>	<b>93.73</b>	<b>93.75</b>
VL2V [1]	92.08	91.90
VL2V+ <b>DIPT</b>	<b>93.46</b>	<b>93.42</b>

**Results.** Table 1 summarizes the performance of baselines trained with DIPT-learned prompts. The original RISE and VL2V baselines are reported using an aggregated template embeddings and a single generic prompt, respectively, following their original implementations. On the Camelyon17-WILDS dataset, our method improves accuracy by up to **6.43%**, while according to Table 2, the worst-case F1 score experiences an improvement of up to **6.05%** (VL2V\*). Also the mean accuracy improves by up to **2.59%**, and the mean F1 score increases by up to **3.08%**, with consistent gains observed across different configurations. Similarly, on the Kather19 dataset, Table 3 shows that combining our DIPT approach with KD methods further improves the test set F1 score by **1.52%**.

## 5 Conclusions

Our study shows that incorporating DIPT step before KD pipelines enhances domain generalization. DIPT entails learning and aggregating domain-specific

prompts in the embedding space. This aids in generation of a class-generic representation. Our approach improves model generalization and handles domain shifts in the histopathology, where domain definitions do not exist.

## References

1. Addepalli, S., Asokan, A.R., Sharma, L., Babu, R.V.: Leveraging vision-language models for improving domain generalization in image classification. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 23922–23932 (2024)
2. Chen, H., Wang, Y., Xu, C., Xu, C., Tao, D.: Learning student networks via feature embedding. *IEEE Transactions on Neural Networks and Learning Systems* **32**(1), 25–35 (2020)
3. Fang, A., Ilharco, G., Wortsman, M., Wan, Y., Shankar, V., Dave, A., Schmidt, L.: Data determines distributional robustness in contrastive language image pre-training (clip). In: International Conference on Machine Learning. pp. 6216–6234. PMLR (2022)
4. Hinton, G., Vinyals, O., Dean, J.: Distilling the knowledge in a neural network. arXiv preprint arXiv:1503.02531 (2015)
5. Huang, Z., Zhou, A., Ling, Z., Cai, M., Wang, H., Lee, Y.J.: A sentence speaks a thousand images: Domain generalization through distilling clip with language guidance. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 11685–11695 (2023)
6. Huang, Z., Bianchi, F., Yuksekgonul, M., Montine, T.J., Zou, J.: A visual–language foundation model for pathology image analysis using medical twitter. *Nature medicine* **29**(9), 2307–2316 (2023)
7. Jahanifar, M., Raza, M., Xu, K., Vuong, T., Jewsbury, R., Shephard, A., Zamani-tajeddin, N., Kwak, J.T., Raza, S.E.A., Minhas, F., et al.: Domain generalization in computational pathology: Survey and guidelines. arXiv preprint arXiv:2310.19656 (2023)
8. Kather, J.N., Krisam, J., Charoentong, P., Luedde, T., Herpel, E., Weis, C.A., Gaiser, T., Marx, A., Valous, N.A., Ferber, D., et al.: Predicting survival from colorectal cancer histology slides using deep learning: A retrospective multicenter study. *PLoS medicine* **16**(1), e1002730 (2019)
9. Koh, P.W., Sagawa, S., Marklund, H., Xie, S.M., Zhang, M., Balsubramani, A., Hu, W., Yasunaga, M., Phillips, R.L., Gao, I., et al.: Wilds: A benchmark of in-the-wild distribution shifts. In: International conference on machine learning. pp. 5637–5664. PMLR (2021)
10. Van der Laak, J., Litjens, G., Ciompi, F.: Deep learning in histopathology: the path to the clinic. *Nature medicine* **27**(5), 775–784 (2021)
11. Luo, P., Zhu, Z., Liu, Z., Wang, X., Tang, X.: Face model compression by distilling knowledge from neurons. In: Proceedings of the AAAI conference on artificial intelligence. vol. 30 (2016)
12. Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al.: Learning transferable visual models from natural language supervision. In: International conference on machine learning. pp. 8748–8763. PmLR (2021)

13. Shakeri, F., Huang, Y., Silva-Rodríguez, J., Bahig, H., Tang, A., Dolz, J., Ben Ayed, I.: Few-shot adaptation of medical vision-language models. In: International Conference on Medical Image Computing and Computer-Assisted Intervention. pp. 553–563. Springer (2024)
14. Stacke, K., Eilertsen, G., Unger, J., Lundström, C.: Measuring domain shift for deep learning in histopathology. *IEEE journal of biomedical and health informatics* **25**(2), 325–336 (2020)
15. Wang, Y., Li, H., Chau, L.p., Kot, A.C.: Embracing the dark knowledge: Domain generalization using regularized knowledge distillation. In: Proceedings of the 29th ACM international conference on multimedia. pp. 2595–2604 (2021)
16. Yao, H., Zhang, R., Xu, C.: Visual-language prompt tuning with knowledge-guided context optimization. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 6757–6767 (2023)
17. Zang, Y., Li, W., Zhou, K., Huang, C., Loy, C.C.: Unified vision and language prompt learning. *arXiv preprint arXiv:2210.07225* (2022)
18. Zhou, K., Yang, J., Loy, C.C., Liu, Z.: Conditional prompt learning for vision-language models. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 16816–16825 (2022)
19. Zhou, K., Yang, J., Loy, C.C., Liu, Z.: Learning to prompt for vision-language models. *International Journal of Computer Vision* **130**(9), 2337–2348 (2022)