

Memory-Efficient Prompt Tuning for Incremental Histopathology Classification

Yu Zhu^{1,3*}, Kang Li^{1*†}, Lequan Yu², Pheng-Ann Heng¹,

¹Department of Computer Science and Engineering, The Chinese University of Hong Kong

²Department of Statistics and Actuarial Science, The University of Hong Kong

³Department of Mechanical Engineering, The University of Hong Kong

{yzhu, kli, pheng}@cse.cuhk.edu.hk, lqyu@hku.hk

Abstract

Recent studies have made remarkable progress in histopathology classification. Based on current successes, contemporary works proposed to further upgrade the model towards a more generalizable and robust direction through incrementally learning from the sequentially delivered domains. Unlike previous parameter isolation based approaches that usually demand massive computation resources during model updating, we present a memory-efficient prompt tuning framework to cultivate model generalization potential in economical memory cost. For each incoming domain, we reuse the existing parameters of the initial classification model and attach lightweight trainable prompts into it for customized tuning. Considering the domain heterogeneity, we perform decoupled prompt tuning, where we adopt a domain-specific prompt for each domain to independently investigate its distinctive characteristics, and one domain-invariant prompt shared across all domains to continually explore the common content embedding throughout time. All domain-specific prompts will be appended to the prompt bank and isolated from further changes to prevent forgetting the distinctive features of early-seen domains. While the domain-invariant prompt will be passed on and iteratively evolve by style-augmented prompt refining to improve model generalization capability over time. In specific, we construct a graph with existing prompts and build a style-augmented graph attention network to guide the domain-invariant prompt exploring the overlapped latent embedding among all delivered domains for more domain-generic representations. We have extensively evaluated our framework with two histopathology tasks, *i.e.*, breast cancer metastasis classification and epithelium-stroma tissue classification, where our approach yielded superior performance and memory efficiency over the competing methods.

Introduction

Histopathology classification is a fundamental task in cancer diagnosis. It aims to specify the malignancy and benignity of suspected tissues by microscopic examination. The resulting analysis is normally considered the gold standard in determining the presence and spread of certain cancers (Benjordi et al. 2017). Although recent deep-learning models

have achieved remarkable progress on this task, contemporary studies are not content with the achievements made so far but strive to upgrade and update model functionality toward perfection by incremental learning (Derakhshani et al. 2022; Li, Yu, and Heng 2022).

One practical yet challenging direction for model upgrading is to incrementally boost its generalization potential over heterogeneous histopathology data. Depending on the technician skills and digital scanner brands in different medical centers, the histology data sampled from multiple sites (*i.e.*, domain) often exhibit heterogenous appearances after hematoxylin and eosin (H&E) staining, varying from dark blueish purple to light pinkish purple (Lin et al. 2019). Then, domain incremental learning (DIL), *i.e.*, a model updating paradigm that enables the model to progressively adapt to more and more heterogenous domains as time goes by, would be substantial for robust histopathology classification.

For any updated model, the basic requirement is to keep the existing capability unaffected, *i.e.*, not catastrophically forgetting the previously-acquired domains. Moreover, we expect to enhance its generalization ability, *i.e.*, not only well adapted to the currently delivered domains but also the unseen domains that might be encountered in the future. Particularly in the medical field, for each update, the model would have no access to the early-delivered domains due to data privacy concerns (Li et al. 2020) and storage burden (Lin et al. 2019). In addition, the domain identity, *i.e.*, the label indicating which domain one particular sample comes from, is erased as part of patient privacy during data anonymization and would be unavailable to use for model training and testing during the entire learning lifespan (Gonzalez, Sakas, and Mukhopadhyay 2020).

The straightforward approach is to finetune the previous model with each sequentially incoming domain one by one. However, with the absence of past domains and data heterogeneity throughout time, it inevitably overrides and disrupts the parameters learned for past domains, leading to catastrophic forgetting of them (Li and Hoiem 2017). A promising way for this issue is to address it from the model-centric perspective, *i.e.*, isolating the early-acquired parameters (*e.g.*, the whole model) into separate storage and allocating new parameters to acquire the newly-arrived domain (Gonzalez, Sakas, and Mukhopadhyay 2020; Miao et al. 2022; Li et al. 2019). Despite their effectiveness, most

*These authors contributed equally.

†Corresponding Author.

Copyright © 2024, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

of them are extremely memory-intensive with increasing computation demands and memory usage over time, greatly limiting their applicability in gigapixel-sized histopathology images (*e.g.*, 1-3 GB per slide (Zhao et al. 2019)).

Fortunately, we could borrow some insights from recent advances in prompt-based natural language processing approaches (Su et al. 2022), which pointed out that employing learnable prompt tokens as the parameterized inputs, could encode necessary guidance to conditionally adapt the frozen pre-trained model to the downstream target task. Inspired by that, it would be unnecessary to completely adjust the previous model to accommodate the currently delivered domain, nor isolate the entire model parameters in separate memory units to retain early-acquired knowledge. Alternatively, it could be more memory-efficient to simply perform prompt tuning upon the initial well-trained classification model and save these lightweight prompts for future usage instead.

In this paper, we present a memory-efficient prompt tuning framework to incrementally learn from the sequentially-delivered heterogeneous domains, progressively cultivating the histopathology classification model towards a more generalizable and robust direction over time. Considering the data heterogeneity of the domains delivered in different time steps, we perform decoupled prompt tuning with two types of prompts. We employ a domain-specific prompt for each domain to independently investigate its distinctive features while maintaining a domain-invariant prompt shared across all domains to continually explore the common content embedding over time. For each incoming domain, we freeze the initial model and train two lightweight prompts upon the existing weights for memory and computation efficiency. We learn a domain-specific prompt from scratch and learn the shared domain-invariant prompt iteratively upon the previous one via style-augmented prompt refining. Specifically, we build up a graph with the existing prompts and constrain the domain-invariant prompt to explore the co-existing and domain-agnostic representations among all seen domains via graph attention propagation. Meanwhile, we augment the style variations met in the prompt refining process to expose more domain-generic representations and further boost its generalization potential. At the end of each time step, we store all prompts in the bank under economical memory costs. The domain-specific prompts would be isolated from further changes and retrieved later to prevent forgetting early-seen domains. The domain-invariant prompt would be carried forward to incrementally acquire more domain-generic features to improve generalization ability. We have extensively evaluated our framework on two histopathology classification tasks, including breast cancer metastase classification on the Camelyon17 dataset (Bandi et al. 2018) and epithelium-stroma tissue classification on a multi-site data collection. In both tasks, our approach showed superior performance over competing methods with better generalization on unseen domains and less forgetting of past domains. Our main contributions could be summarized as follows:

- We proposed a memory-efficient prompt tuning framework to iteratively upgrade the model towards a more generalized direction in economical memory cost.
- We performed decoupled prompt tuning with a series of domain-specific prompts and a shared domain-invariant prompt to tackle the heterogeneity of incoming domains.
- We presented style-augmented prompt refining to iteratively evolve the domain-invariant prompt over time to boost its generalization potential on unseen data.
- We have validated our approach on two histopathology image classification tasks, where our framework outperformed other comparison methods significantly.

Related Work

Domain Incremental Learning

Considerable efforts have been devoted to domain incremental learning to progressively cultivate the model accommodating more and more heterogeneous domains. One stream of works would not require any additional module to support model updating (Aljundi et al. 2018; Li and Hoiem 2017; Kirkpatrick et al. 2017; Zenke, Poole, and Ganguli 2017). For example, the regularization-based methods (Kirkpatrick et al. 2017; Aljundi et al. 2018) employed a loss term to penalize large changes of the parameters important to historical domains to help retain early-acquired knowledge. However, these approaches often suffered from interval forgetting when dealing with a long sequence of incremental learning tasks, and their performance still has certain improvement spaces (Luo et al. 2020; Mai et al. 2022). Other streams of work (*e.g.*, replayed-based methods and parameter isolation methods) sacrificed memory usage to trade for better model performance. For example, Shin *et al.* (Shin et al. 2017) employed an extra generative adversarial network (GAN) ($\approx 266\text{MB}$) to memorize and replay past domain distributions to prevent forgetting past domains, while Gonzalez *et al.* (Gonzalez, Sakas, and Mukhopadhyay 2020) stored all previously-learned models ($\approx 81\text{MB}$ each) in a separate space and maintain an autoencoder-based domain classifier to retrieve them back when necessary. Although the above methods could effectively alleviate the forgetting of historical domains, it also results in massive memory consumption, making them less applicable to gigabyte-size histopathology images. In contrast, we maximally reuse the existing initial model and perform decoupled prompt tuning upon it by two lightweight prompts ($\approx 0.5\text{MB}$) for each incoming domain, greatly boosting memory efficiency.

Prompt Learning

Inspired by the recent progress of prompt tuning in natural language processing (Su et al. 2022), contemporary works attempted to apply it for incremental learning. Most prior works concentrated on the class incremental learning (CIL) settings, *i.e.*, progressively learning to categorize more and more classes over time. They prevented forgetting early-acquired classes by creating a shared prompt pool for instance-wise prompt query (Wang et al. 2022b), setting general prompts and expert prompts to form complementary learning (Wang et al. 2022a) and *etc.* However, most of them are less prepared for the domain incremental learning settings, especially the demand to generalize to unseen

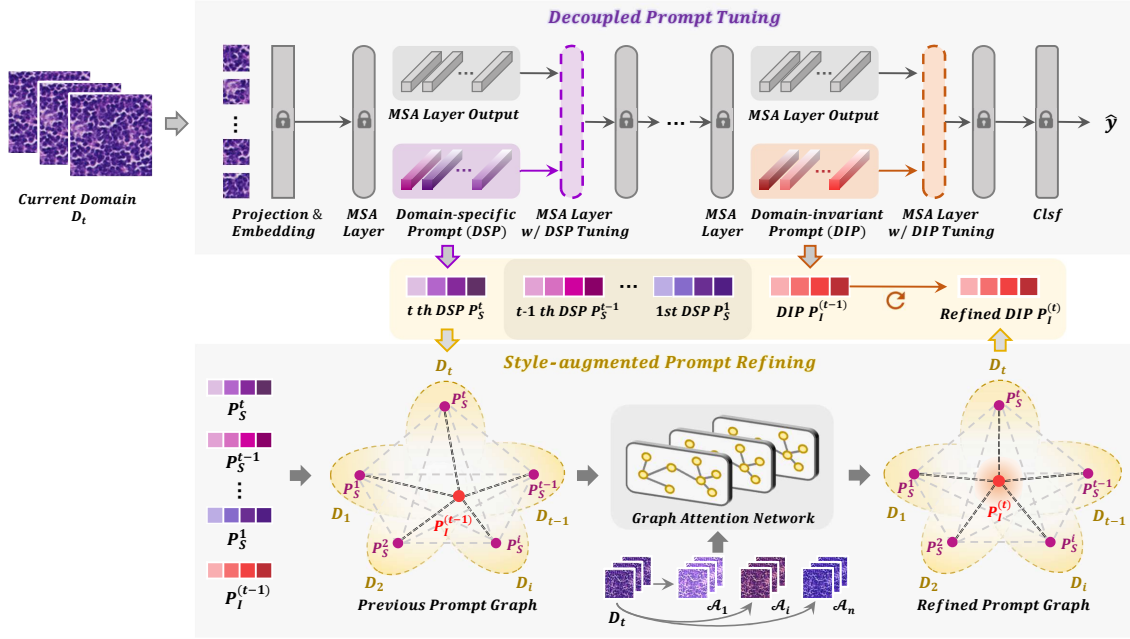


Figure 1: Overview of our memory-efficient prompt tuning framework. We proposed to perform decoupled prompt tuning upon the initial model with two lightweight prompts, aiming to acquire the latest domain knowledge in economical memory cost. We employ a domain-specific prompt independently for each domain to acquire its distinctive features like appearances. The learned domain-specific prompt would be stored and isolated in the prompt bank to help alleviate the forgetting of early-acquired domains. Meanwhile, we maintain a domain-invariant prompt shared across domains to progressively learn the common content over time like shape prior. We performed style-augmented prompt refining upon the previous domain-invariant prompt, where we constrain its exploration scope within the overlapped latent embeddings of all seen domains and guide it to learn the domain-generic representations, gradually strengthening the generalization potential over time.

data. Most CIL approaches would not expect the model to correctly recognize the objects of unseen classes (*i.e.*, never learned in training). However, in domain incremental learning settings, it is highly desired for a model well generalizing to unseen domains of unknown appearances for robust classification. Very recently, a DIL approach S-Prompt (Wang, Huang, and Hong 2022) tried to independently learn the prompts across domains for a win-win game but still overlooked the generalization issue. In our work, we put extra effort to maintain a domain-variant prompt shared over time by style-augmented prompt refining, incrementally absorbing more domain-generic features to improve generalization.

Methodology

In domain incremental learning (DIL) settings, we assume a heterogenous data stream D_1, D_2, \dots, D_T sequentially delivered from multiple sites one by one. With the arrival of the dataset D_t at time step t , our goal is to incrementally optimize the previous model M_{t-1} with D_t , such that the updated model M_t would not catastrophically forget past domains D_1, D_2, \dots, D_{t-1} , while maintaining satisfying generalization ability for unseen domains. For privacy concerns in medical fields, all past domains would be inaccessible and no domain identity would be available.

Fig. 1 overviews our framework. For each incoming domain, we reuse the initial model and perform decoupled

prompt tuning upon it with two lightweight prompts to acquire new domain knowledge in a memory-efficient manner. In specific, the domain-specific prompt (DSP) is independently learned from scratch to tackle the distinctive features, while the domain-invariant prompt (DIP) is iteratively evolved from the previous one by style-augmented prompt refining to incrementally explore domain-generic features.

Decoupled Prompt Tuning

We construct a transformer backbone (*e.g.*, ViT (Dosovitskiy et al. 2020)) for the classification model. It consists of a basic transformer feature extractor f_b to convert the input image into sequence-like high-level representations, and a classification layer f_ϕ to map the representation to the final prediction \hat{y} . At time step t , with the arrival of the current domain D_t , we load the pre-trained weights into the basic feature extractor following prior works (Wang et al. 2022a,b) and freeze them. Upon it, we perform decoupled prompt tuning by two lightweight trainable prompts, *i.e.*, one domain-invariant prompt $p_I^{(t)}$ and one domain-specific prompt p_s^t , to acquire the current domain. To avoid any confusion, we use the superscript (t) to denote the shared domain-invariant prompt learned in the t -th time step, while using the superscript t to indicate the t -th domain-specific prompt.

The domain-invariant prompt and the domain-specific prompt can be inserted as additional inputs of any multi-

head self-attention (MSA) layer in the basic transformer feature extractor. Take the i -th MSA layer as an example. Before passing the previous MSA layer outputs $h_{i-1} \in \mathbb{R}^{l \times m}$ to it, we keep the query h_{i-1}^q and append the domain-invariant prompt $p_I^{(t)} \in \mathbb{R}^{l \times m}$ in its key h_{i-1}^k and value h_{i-1}^v to guide it explore the domain-shared representations as

$$h_i = f_{MSA}^{(i)}(h_{i-1}^q, [p_I^k; h_{i-1}^k], [p_I^v; h_{i-1}^v]), \quad (1)$$

where $f_{MSA}^{(i)}$ and h_i denote the i -th MSA layer and its output tuned with the domain-invariant prompt respectively. $p_I^k \in \mathbb{R}^{l/2 \times m}$ and $p_I^v \in \mathbb{R}^{l/2 \times m}$ are split from $p_I^{(t)}$ to maintain the same sequence length before and after the MSA layer. The domain-specific prompt $p_s^t \in \mathbb{R}^{l \times m}$ can be attached in a similar way to learn the distinctive features as

$$h_j = f_{MSA}^{(j)}(h_{j-1}^q, [p_s^k; h_{j-1}^k], [p_s^v; h_{j-1}^v]), \quad (2)$$

where $f_{MSA}^{(j)}$ and h_j denote the j -th MSA layer and its outputs respectively. $p_s^k, p_s^v \in \mathbb{R}^{l/2 \times m}$ are split from p_s^t .

As domain identity is not available during inference, we additionally equip a distinguishable key value k_t for the domain-specific prompt, to help pair each test image with a matching domain-specific prompt. We decompose each image $x_i \in D_t$ into the amplitude spectrum $\mathcal{A}(x_i)$ and phase spectrum $\mathcal{C}(x_i)$ in the frequency space by fast Fourier transform Φ_{FFT} . Since the amplitude captures the low-level statistics (*e.g.*, style, appearances) while the phase extracts the high-level features (*e.g.*, content, shape) (Jiang, Wang, and Dou 2022; Liu et al. 2021), we implement the key value k_t as the average amplitude spectrum of the images in D_t as

$$k_t = \frac{1}{N_t} \sum_{i=1}^{N_t} \mathcal{A}(x_i), \quad (3)$$

where N_t denotes the total number of training data in D_t .

In the first time step ($t = 1$), we simultaneously optimize the domain-specific prompt p_s^1 , the domain-invariant prompt $p_I^{(1)}$ and the classification layer f_ϕ upon the frozen basic feature extractor f_b by the training samples $(x, y) \in D_1$ as

$$\min_{f_\phi, p_I^{(1)}, p_s^1} \mathcal{L}_{ce} \left(f_\phi \left(f_b \left(x; p_I^{(1)}, p_s^1 \right) \right), y \right), \quad (4)$$

where \mathcal{L}_{ce} denotes the cross-entropy loss. Before moving to the next time step, we store all prompts and the associated keys into the prompt bank $\mathcal{P}^1 = \{p_I^{(1)}, [p_s^1, k_1]\}$. The domain-specific prompt and its key value would be isolated from further changes while the domain-invariant prompt would be passed to the next time step to iteratively evolve.

For the subsequent time step t ($t > 1$), we keep the classification model (including f_ϕ and f_b) frozen, and optimize the domain-specific prompt p_s^t and the domain-invariant prompt $p_I^{(t)}$ asynchronously in separate steps. We first learn an independent domain-specific prompt p_s^t from scratch with the old domain-invariant prompt $p_I^{(t-1)}$ fixed by

$$\min_{p_s^t} \mathcal{L}_{ce} \left(f_\phi \left(f_b \left(x; p_I^{(t-1)}, p_s^t \right) \right), y \right), \quad (5)$$

where $(x, y) \in D_t$. Then we update the domain-invariant prompt by style-augmented prompt refining, which would be thoroughly described in the following subsection.

Style-augmented Prompt Refining

The straightforward way to update the domain-invariant prompt is to finetune the previous one $p_I^{(t-1)}$ along with the t -th domain-specific prompt p_s^t . However, this would easily make the latest domain-invariant prompt not compatible with early domain-specific prompts recorded in the prompt bank. To address this issue, we build a graph with all existing prompts and feed it into the graph attention network (GAT) (Veličković et al. 2017) to guide the domain-invariant prompt exploring the co-existing and generic features.

GAT setup We flatten all existing prompts into long vectors as $P^t = \{p_I^{(t-1)}, p_s^1, \dots, p_s^t\}$, and take them as the nodes of the graph. The graph attention network consists of one learnable linear transformation $\mathbf{W} \in \mathbb{R}^{L \times L}$, where $L = l \times m$, and a trainable single-layer feed-forward neural network a to acquire the attention coefficients e between two nodes. Particularly for the node of the domain-invariant prompt that we most concern, the attention coefficients for its i -th neighbor e_{IS}^i and itself e_{II} are computed as

$$\begin{aligned} e_{IS}^i &= a \left(\mathbf{W} p_I^{(t-1)}, \mathbf{W} p_s^i \right), \\ e_{II} &= a \left(\mathbf{W} p_I^{(t-1)}, \mathbf{W} p_I^{(t-1)} \right), \end{aligned} \quad (6)$$

which indicates the correlation and relevance between the domain-invariant prompt and each prompt in the bank. We normalize the above coefficients as α_{IS}^i and α_{II} , and use them to adjust the participation of each node in the knowledge aggregation to the domain-invariant prompt as

$$\begin{aligned} \alpha_{IS}^i &= \frac{\exp(e_{IS}^i)}{\exp(e_{II}) + \sum_{i=1}^t \exp(e_{IS}^i)}, \\ \alpha_{II} &= \frac{\exp(e_{II})}{\exp(e_{II}) + \sum_{i=1}^t \exp(e_{IS}^i)}, \end{aligned} \quad (7)$$

The outputs of the domain-invariant prompt $p_I^{(t)}$ would be

$$p_I^{(t)} = f_{GAT}(P^t) = \sum_{i=1}^t \alpha_{IS}^i \mathbf{W} p_s^i + \alpha_{II} \mathbf{W} p_I^{(t-1)}, \quad (8)$$

where f_{GAT} denotes the graph attention network. By simply reshaping $p_I^{(t)}$ into the original prompt size, we could obtain the updated domain-invariant prompt $p_I^{(t)}$.

Style-augmented GAT training We augment the style diversity met in GAT training to further improve the generalization potential of the domain-invariant prompt. As aforementioned, for each image, its amplitude spectrum reflects the low-level features like style or appearance, while the phase spectrum presents its high-level content like shape. Given an image-label pair of the current domain $(x_i, y_i) \in D_t$, we reserve the phrase spectrum $\mathcal{C}(x_i)$ to keep its semantic content, but substitute its amplitude $\mathcal{A}(x_i)$ into a new one $\mathcal{A}'(x_i)$ to modulate its style. As shown in Fig. 2, rather than blindly guessing feasible appearances, we use a set of random scalars $\{\lambda_i^1, \dots, \lambda_i^t\}$ to interpolate the average amplitudes of past domains (*i.e.*, the keys in the prompt bank)

and generate the new amplitude $\mathcal{A}'(x_i)$ as

$$\mathcal{A}'(x_i) = \lambda_i^1 k_1 + \lambda_i^2 k_2 + \dots + \lambda_i^t k_t, \quad (9)$$

where $\sum_{j=0}^t \lambda_i^j = 1$. We then perform inverse fast Fourier transform operation Φ_{FFT}^{-1} to remap the phase $\mathcal{C}(x_i)$ and amplitude $\mathcal{A}'(x_i)$ into the image space, and generate the style-augmented image x'_i as follows

$$x'_i = \Phi_{\text{FFT}}^{-1}(\mathcal{A}'(x_i), \mathcal{C}(x_i)). \quad (10)$$

We paired the style-augmented image x'_i with its original label y_i to form a set $D_t^{s^a} = \{(x'_i, y_i), i \in [1, N_t]\}$, which will be used for GAT training along with the current domain D_t . For any data $(x, y) \in D_t \cup D_t^{s^a}$, we select the most compatible domain-specific prompt $p_s^{j^*}$ by similarity ranking as

$$j^* = \operatorname{argmax}_j \gamma(\mathcal{A}(x), k_j), \quad (11)$$

where $j \in [1, t]$ and γ denotes the cosine similarity. Then, we force the GAT to produce a domain-invariant prompt p_I^t that could satisfyingly tackle the images of any augmented style and work smoothly with any domain-specific prompts in the bank by the following objectives

$$\min_{f_{\text{GAT}}} \mathcal{L}_{ce}(f_\phi(f_b(x; f_{\text{GAT}}(P^t), p_s^{j^*})), y). \quad (12)$$

The overall training scheme is presented in Algorithm 1. During inference, we use the latest domain-invariant prompt for all test data and pair each test sample with the most compatible domain-specific prompt to it by Eq. 11 accordingly.

Experiment

Dataset and Experiment Settings

Breast cancer metastase classification We adopted the Camelyon17 dataset (Bandi et al. 2018) which provided the labels of the presence or absence of breast cancer. The data was collected from 5 medical centers with different stains. We closely followed the domain split of prior works (Jiang, Wang, and Dou 2022) and took the samples of the same center as one domain. All domains are sequentially delivered one by one in ascending order. We set the total time step as 4, where Domain 4 currently arrives, Domain 1-3 are previously delivered, and Domain 5 remains unseen to the model.

Epithelium-stroma tissue classification We utilized four public datasets, including 615 images from VGH (Beck et al. 2011) (Domain 1), 671 images from NKI (Beck et al. 2011) (Domain 2), 1296 patches from IHC (Linder et al. 2012) (Domain 3), and 26,437 patches from NCH (Kather et al. 2019) (Domain 4). Each of them comes from different institutions under different H&E stains. Here, we set the total time step as 3, where Domain 3 has currently arrived, Domain 1 and 2 are previously delivered and Domain 4 remains unseen during model training.

Implementation details We adopted the ViT-B/16 (Dosovitskiy et al. 2020) as our feature extractor f_b . We employ the Adam optimizer with the learning rate of $7.5e^{-4}$ in the first time step and the learning rate of $1e^{-4}$ for the subsequent time steps.

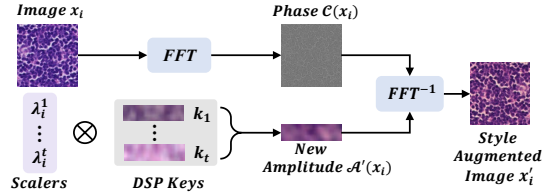


Figure 2: Illustration of generating style-augmented data.

Algorithm 1: Training Procedures

Output: The model $f_\phi(f_b(\cdot))$ and prompt bank \mathcal{P} .

while incrementally learning from $t = 1$ to T **do**

if $t == 1$ **then**

 Load pre-trained weights in f_b and freeze it.

 Optimize $p_I^{(1)}, p_s^1, f_\phi$ with D_1 by Eq. 4.

 Calculate the key k_1 for p_s^1 by Eq. 3.

 Store all prompts in the bank \mathcal{P}^1 .

else

 Freeze f_b, f_ϕ and train p_s^t with D_t by Eq. 5.

 Generate style-augmented data $D_t^{s^a}$ by Eq. 9.

 Update $p_I^{(t)}$ given p_s^t by GAT as Eq. 12.

 Compute the key k_t by Eq. 9 and append $[k_t, p_s^t]$ in the prompt bank \mathcal{P}^t .

 Overwrite DIP as $p_I^{(t)}$ in the prompt bank \mathcal{P}^t .

end

 Pass f_b, f_ϕ and \mathcal{P}^t in the $t + 1$ step.

end

Return $f_\phi(f_b(\cdot)), \mathcal{P} \leftarrow \mathcal{P}^T$.

Evaluation metrics We employed the classification accuracy (Acc) as the base evaluation metric. We first measure the model performance at the last incremental learning step on all domains, including previous domains, the current domain and unseen domains, to extensively evaluate the ability to alleviate forgetting and generalize. We further employ three more metrics to comprehensively evaluate the overall performance of the entire incremental learning span. Backward transfer (BWT) evaluates the model stability, *i.e.*, the ability to alleviate catastrophic forgetting, which is computed as $BWT = 2 \sum_{i=2}^N \sum_{j=1}^{i-1} (R_{i,j} - R_{j,j}) / N(N-1)$, where $R_{i,j}$ denotes the Acc of the model trained sequentially from the 1st domain to the i -th domain and tested on the j -th domain, and N denotes the total number of training domains. Incremental learning (IL) provides an overall measurement of all seen domains and measures both stability and plasticity of the model as $IL = 2 \sum_{i=1}^N \sum_{j=1}^i R_{i,j} / N(N+1)$. Forward Transfer on unseen domains (FTU) measures the generalization ability of the model as $FTU = 2 \sum_{i=1}^N \sum_{j=i+1}^N R_{i,j} / N(N-1)$.

Comparison with the State-of-the-arts

Compared methods Besides the intuitive method, *i.e.*, sequential finetune, we also implemented the state-of-the-art incremental learning methods, including the regularization-based methods (LwF (Li and Hoiem 2017), EWC (Kirk-

Methods	Acc [%] ↑						IL ↑	BWT ↑	FTU ↑
	Domain 1	Domain 2	Domain 3	Domain 4	Domain 5	Avg			
Individual Training	96.79 (±0.77)	91.62 (±0.82)	96.32 (±0.54)	96.46 (±0.65)	96.84 (±0.70)	95.30 (±0.52)	71.74 (±0.56)	-38.65 (±0.55)	64.71 (±1.31)
Joint Training (Upper bound)	96.96 (±0.54)	94.32 (±0.67)	97.51 (±0.52)	97.41 (±0.63)	82.80 (±1.19)	92.45 (±0.58)	95.27 (±0.46)	2.53 (±0.03)	80.75 (±1.22)
Sequential Finetune	61.92 (±1.22)	64.78 (±2.55)	51.52 (±1.10)	96.51 (±0.68)	49.82 (±3.22)	64.91 (±2.84)	72.59 (±2.01)	-37.13 (±1.69)	63.69 (±2.62)
LwF	83.28 (±0.44)	72.69 (±1.05)	59.71 (±0.65)	95.86 (±0.83)	69.07 (±0.53)	76.12 (±0.68)	80.12 (±1.11)	-23.44 (±0.36)	69.18 (±0.47)
EWC	81.87 (±0.81)	77.21 (±1.44)	53.26 (±2.21)	96.47 (±0.99)	46.77 (±2.55)	71.12 (±2.28)	77.78 (±2.46)	-29.32 (±0.48)	61.02 (±2.26)
SI	86.53 (±1.66)	82.49 (±0.97)	54.04 (±1.78)	96.52 (±0.64)	50.55 (±1.43)	74.03 (±1.21)	79.83 (±1.29)	-24.15 (±0.61)	62.51 (±1.05)
DGR	81.21 (±2.49)	83.20 (±1.96)	72.77 (±3.45)	96.84 (±1.10)	78.57 (±2.83)	82.52 (±2.36)	86.80 (±1.95)	-15.67 (±1.62)	70.64 (±2.18)
Orc-MML	90.07 (±0.93)	86.64 (±0.85)	87.32 (±0.77)	89.73 (±0.96)	66.91 (±1.59)	84.13 (±0.91)	89.69 (±0.94)	-6.44 (±0.60)	75.45 (±1.03)
S-Prompt	92.50 (±0.69)	81.30 (±0.88)	93.93 (±0.70)	94.43 (±0.83)	74.37 (±0.92)	87.31 (±0.85)	91.73 (±0.78)	-3.46 (±0.24)	80.64 (±0.62)
DualPrompt	91.34 (±0.43)	85.18 (±0.75)	92.71 (±0.67)	95.63 (±0.58)	67.44 (±2.78)	86.46 (±0.59)	91.11 (±0.42)	-7.60 (±0.11)	77.57 (±1.02)
Ours	94.36 (±0.90)	89.40 (±0.39)	94.32 (±0.87)	95.86 (±0.79)	84.12 (±1.26)	91.61 (±0.67)	93.67 (±0.44)	-1.62 (±0.09)	82.17 (±0.54)

Table 1: The comparison results of breast cancer metastases classification in the final time step (the 2-7 Columns) and the entire domain incremental learning process (the last three columns). We have highlighted the best DIL results in bold.

patrick et al. 2017) and SI (Zenke, Poole, and Ganguli 2017)), the replay-based method (DGR (Shin et al. 2017)), and the parameter isolation methods (Orc-MML (Gonzalez, Sakas, and Mukhopadhyay 2020)) especially those also involved prompt tuning (DualPrompt (Wang et al. 2022a) and S-Prompt (Wang, Huang, and Hong 2022)). We separately trained a model for each domain (*i.e.*, individual training) and also jointly trained a model with all delivered domains. Since the joint training could fully access all seen domains while the DIL methods only have access to the current domain, we consider it as the upper bound.

Experiment results We reported the performance of breast cancer metastase classification in Table 1. Each experiment is repeated 5 times to avoid random bias. Our approach achieved the best results in the majority of evaluation metrics (8 out of 9). Compared to prior parameter isolation approaches like Orc-MML, DualPrompt, and S-Prompt, our framework yielded the least forgetting of past domains with 1.86% increases in Domain 1 and 2.76% increases in Domain 2. Since we continually evolve the domain-invariant prompt by style-augmented prompt refining, our approach greatly improved model generalization capability, leading to 5.55% gains in the unseen domain (Domain 5) and 1.53% increases in FTU compared to the state-of-the-art approach.

We presented the results of epithelium-stroma tissue classification in Table 2. Our framework outperformed the competing approaches on most of the evaluation metrics (7 out of 8). For the first delivered domain (Domain 1), which

commonly suffered most from catastrophic forgetting, our framework achieved 5.32% higher than S-Prompt and 1.67% higher than DualPrompt, demonstrating the effectiveness of our decoupled prompt tuning. When it comes to the generalization ability, our approach yielded 2.42% increases in the unseen domain (Domain 4) and 1.77% increases in FTU over the competing methods.

Analysis of the Key Components

The tradeoff between model performance and memory efficiency We evaluated the memory efficiency by the metric of model size efficiency (MS) (Díaz-Rodríguez et al. 2018), which measures the additional storage used at the time step t compared to the usage in the first time step by computing $MS = \min\left(1, \frac{1}{N} \sum_{i=1}^N \frac{\theta_t}{\theta_i}\right)$, where θ_1 and θ_i denote the allocated memory spaces to store all necessary modules for the next round of incremental learning in the 1st and i -th time step respectively. We also calculated the absolute value of the average additional memory storage (AAMS) over time as $AAMS = \frac{1}{N} \sum_{i=1}^N |\theta_i - \theta_1|$.

As presented in Table 3, the sequential finetune (Seq-FT) method and regularization-based approaches barely require any additional module to support the learning in the next time step. However, these methods still have large improvement spaces in the aspect of alleviating model forgetting (see Acc and BWT) and enhancing model generalization (see FTU). The replay-based approaches and param-

Methods	Acc [%] \uparrow					IL \uparrow	BWT \uparrow	FTU \uparrow
	Domain 1	Domain 2	Domain 3	Domain 4	Avg			
Individual Training	94.09 (± 0.56)	91.84 (± 0.61)	93.42 (± 0.83)	96.50 (± 0.69)	93.96 (± 0.60)	78.79 (± 0.64)	-28.87 (± 0.35)	67.74 (± 2.19)
Joint Training (Upper bound)	94.33 (± 0.48)	92.61 (± 0.65)	93.67 (± 0.71)	85.82 (± 1.86)	91.61 (± 0.84)	93.56 (± 0.52)	1.20 (± 0.04)	74.91 (± 0.93)
Sequential Finetune	62.77 (± 1.19)	70.14 (± 2.06)	93.58 (± 0.95)	69.06 (± 1.42)	73.89 (± 1.27)	82.84 (± 1.44)	-20.86 (± 1.08)	66.02 (± 1.79)
LwF	77.54 (± 1.60)	66.32 (± 2.47)	93.52 (± 1.01)	73.61 (± 1.99)	77.75 (± 1.32)	85.12 (± 1.17)	-15.82 (± 0.91)	73.93 (± 1.85)
EWC	71.22 (± 2.04)	77.90 (± 1.70)	93.61 (± 0.84)	68.12 (± 2.66)	77.71 (± 1.82)	87.85 (± 1.79)	-15.01 (± 0.38)	72.32 (± 1.88)
SI	66.58 (± 0.47)	72.51 (± 0.21)	93.55 (± 1.02)	77.33 (± 1.55)	77.50 (± 0.68)	82.64 (± 0.64)	-21.99 (± 0.82)	71.51 (± 0.73)
DGR	80.54 (± 1.30)	82.35 (± 1.77)	93.67 (± 0.69)	86.72 (± 1.48)	85.81 (± 0.74)	88.80 (± 0.60)	-9.72 (± 0.18)	78.29 (± 1.95)
Orc-MML	85.61 (± 0.93)	84.93 (± 0.75)	83.76 (± 0.67)	72.58 (± 0.99)	81.72 (± 0.72)	88.15 (± 0.71)	-6.82 (± 0.05)	69.19 (± 1.26)
S-Prompt	85.29 (± 0.83)	88.42 (± 1.94)	94.33 (± 0.75)	73.03 (± 1.76)	85.27 (± 0.94)	89.61 (± 1.04)	-4.19 (± 0.06)	74.62 (± 1.01)
DualPrompt	88.94 (± 0.58)	88.05 (± 0.47)	93.63 (± 0.17)	83.16 (± 0.51)	88.45 (± 0.33)	91.14 (± 0.77)	-4.34 (± 0.03)	75.11 (± 0.42)
Ours	90.61 (± 1.00)	88.47 (± 0.39)	93.84 (± 1.03)	89.14 (± 0.97)	90.52 (± 0.55)	92.17 (± 0.49)	-2.19 (± 0.05)	80.06 (± 0.84)

Table 2: The comparison results of epithelium-stroma classification in the final time step (the 2-6 columns) and the entire domain incremental learning process (the last three columns). We have highlighted the best DIL results in bold.

ter isolation approaches, such as Orc-MML (Orc-M), normally consumed extra memory spaces to trade for model performance. Among them, the prompt-based approaches, *i.e.*, S-Prompt (S-P), DualPrompt (Dual-P), and ours, are the top 3 memory-efficient ones. With limited additional memory spaces (around 0.57 MB), our approach could bring significant performance gains with 4.30% in the average Acc, 1.84% in BWT, and 1.53% in FTU over prior prompt-based approaches, suggesting it is the most desirable approach when considering the trade-off between model accuracy and memory consumption.

For training efficiency, our work used 0.6h longer training time than other prompt-based methods on average, which is generally affordable in most cases.

The efficiency of decoupled prompt tuning We visualized the output feature embeddings after performing decoupled prompt tuning in the breast cancer classification task via t-SNE in Fig. 3. For the embeddings of the same category (*e.g.*, Tumor or Normal), the features within the same domain are grouped into a cluster and well-separated from other domains, suggesting that the learned domain-specific prompts could effectively capture the domain-distinctive characteristics. For the embeddings within the same domain, it shows a clear decision boundary between the features of normal tissues and tumor tissues, indicating our model could well distinguish them from each other.

Methods	A-Acc \uparrow	BWT \uparrow	FTU \uparrow	AAMS \downarrow	MS \uparrow
Seq-FT	64.91	-37.13	63.69	0	1
LwF	76.12	-23.44	69.18	0	1
EWC	71.12	-29.32	61.02	0	1
SI	74.03	-24.15	62.51	0	1
DGR	82.52	-15.67	70.64	399.61	0.56
Orc-M	84.13	-6.44	75.45	460.50	0.52
S-P	87.31	-3.46	80.64	0.23	0.99
Dual-P	86.46	-7.60	77.57	0.57	0.99
Ours	91.61	-1.62	82.17	0.57	0.99

Table 3: Analysis of model accuracy and memory efficiency. Here, we reported the average Acc of all domains (A-Acc) in the last time step, and employed model size efficiency (MS) and the average additional memory storage (AAMS) [MB] to measure memory efficiency.

The study of key operations in style-augmented prompt refining To extensively investigate the effectiveness of style-augmented prompt refining, we experimented on several settings, including (a) using DIP individually (first row), (b) using DSP individually (second row), (c) using both DIP and DSP but updating DIP via fine-tuning (third row), and (d) using both DIP and DSP and refining DIP via GAT with the current domain D_t only (fourth row), and compared them with ours, *i.e.*, using both DIP and DSP and refining

DIP	DSP	GAT	SA	A-Acc \uparrow	BWT \uparrow	FTU \uparrow
✓				85.94	-8.86	77.35
	✓			86.31	-7.29	75.82
✓	✓			86.24	-7.16	77.49
✓	✓	✓		90.34	-3.19	80.04
✓	✓	✓	✓	91.61	-1.62	82.17

Table 4: Analysis of the key operations in style-augmented prompt refining.

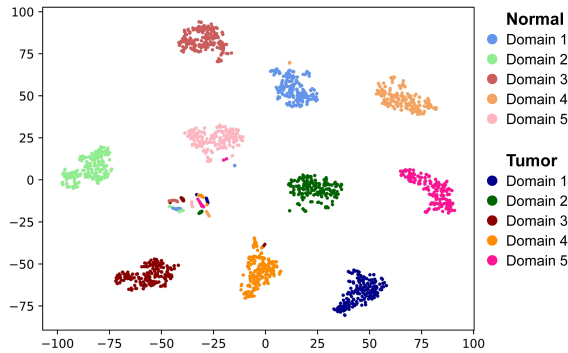


Figure 3: The t-SNE visualization of the feature embeddings after applying decoupled prompt tuning.

DIP via GAT with style-augmented training data $D_t^{sa} \cup D_t$ (the last row). We reported the results on breast cancer classification in Table 4. Compared to simply finetuning DIP, refining with GAT could explore more high-correlative and domain-generic representations across domains, leading to 4.10% increases in the average Acc and 2.55% increases in FTU. Further refining by the style-augmented data not only keeps the updated DIP compatible with early-recorded DSPs but also lets the model be early-prepared for the unseen styles during inference, thus bringing further improvements of 1.27% and 2.13% increases in the average Acc and FTU respectively.

Conclusion

We presented a memory-efficient prompt tuning framework to incrementally evolve the histology classification model towards a more generalizable and robust direction. For each incoming domain, we performed decoupled prompt tuning upon the initial classification model with two lightweight prompts, efficiently acquiring the latest domain knowledge without huge memory costs. We customized a domain-specific prompt customized for tackling the distinctive characteristics while maintaining a domain-invariant prompt shared across all domains to progressively explore the common content embedding. We additionally conducted style-augmented prompt refining on the domain-invariant prompt to continually investigate domain-generic representations across domains and cultivate its generalization potential. All prompts will be stored in a prompt bank, where the domain-specific prompts will be isolated from further changes to prevent catastrophic forgetting of past domains, while the

domain-invariant prompt will be passed on to the next time step to continually evolve. We have extensively evaluated our framework with two histology classification tasks, where our approach outperformed other comparison methods with higher accuracy and more satisfying memory efficiency.

Acknowledgements

The work described in this paper was supported in part by the following grant from the Research Grants Council of the Hong Kong SAR, China (Project No. T45-401/22-N), the Hong Kong Innovation and Technology Fund (Project No. MHP/085/21), and the National Natural Science Fund (62201483).

References

- Aljundi, R.; Babiloni, F.; Elhoseiny, M.; Rohrbach, M.; and Tuytelaars, T. 2018. Memory aware synapses: Learning what (not) to forget. In *ECCV*, 139–154.
- Bandi, P.; Geessink, O.; Manson, Q.; Van Dijk, M.; Balkenhol, M.; Hermsen, M.; et al. 2018. From detection of individual metastases to classification of lymph node status at the patient level: the camelyon17 challenge. *IEEE Transactions on medical imaging*, 38(2): 550–560.
- Beck, A. H.; Sangoi, A. R.; Leung, S.; Marinelli, R. J.; Nielsen, T. O.; Van De Vijver, M. J.; et al. 2011. Systematic analysis of breast cancer morphology uncovers stromal features associated with survival. *Science translational medicine*, 3(108): 108ra113–108ra113.
- Bejnordi, B. E.; Veta, M.; Van Diest, P. J.; Van Ginneken, B.; Karssemeijer, N.; Litjens, G.; et al. 2017. Diagnostic assessment of deep learning algorithms for detection of lymph node metastases in women with breast cancer. *JAMA*, 318(22): 2199–2210.
- Derakhshani, M. M.; Najdenkoska, I.; van Sonsbeek, T.; Zhen, X.; Mahapatra, D.; Worring, M.; et al. 2022. Life-Longer: A Benchmark for Continual Disease Classification. In *MICCAI*, 314–324. Springer.
- Díaz-Rodríguez, N.; Lomonaco, V.; Filliat, D.; and Maltoni, D. 2018. Don’t forget, there is more than forgetting: new metrics for Continual Learning. *arXiv preprint arXiv:1810.13166*.
- Dosovitskiy, A.; Beyer, L.; Kolesnikov, A.; Weissenborn, D.; Zhai, X.; Unterthiner, T.; et al. 2020. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*.
- Gonzalez, C.; Sakas, G.; and Mukhopadhyay, A. 2020. What is Wrong with Continual Learning in Medical Image Segmentation? *arXiv preprint arXiv:2010.11008*.
- Jiang, M.; Wang, Z.; and Dou, Q. 2022. Harmoff: Harmonizing local and global drifts in federated learning on heterogeneous medical images. In *AAAI*, volume 36, 1087–1095.
- Kather, J. N.; Krisam, J.; Charoentong, P.; Luedde, T.; Herpel, E.; Weis, C.-A.; et al. 2019. Predicting survival from colorectal cancer histology slides using deep learning: A retrospective multicenter study. *PLoS medicine*, 16(1): e1002730.

- Kirkpatrick, J.; Pascanu, R.; Rabinowitz, N.; Veness, J.; Desjardins, G.; Rusu, A. A.; et al. 2017. Overcoming catastrophic forgetting in neural networks. *PNAS*, 114(13): 3521–3526.
- Li, K.; Yu, L.; and Heng, P.-A. 2022. Domain-incremental Cardiac Image Segmentation with Style-oriented Replay and Domain-sensitive Feature Whitening. *IEEE Transactions on Medical Imaging*.
- Li, X.; Jiang, M.; Zhang, X.; Kamp, M.; and Dou, Q. 2020. FedBN: Federated Learning on Non-IID Features via Local Batch Normalization. In *ICLR*.
- Li, X.; Zhou, Y.; Wu, T.; Socher, R.; and Xiong, C. 2019. Learn to grow: A continual structure learning framework for overcoming catastrophic forgetting. In *ICML*, 3925–3934. PMLR.
- Li, Z.; and Hoiem, D. 2017. Learning without forgetting. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 40(12): 2935–2947.
- Lin, H.; Chen, H.; Graham, S.; Dou, Q.; Rajpoot, N.; and Heng, P.-A. 2019. Fast scannet: Fast and dense analysis of multi-gigapixel whole-slide images for cancer metastasis detection. *IEEE Transactions on medical imaging*, 38(8): 1948–1958.
- Linder, N.; Konsti, J.; Turkki, R.; Rahtu, E.; Lundin, M.; Nordling, S.; et al. 2012. Identification of tumor epithelium and stroma in tissue microarrays using texture analysis. *Diagnostic pathology*, 7: 1–11.
- Liu, Q.; Chen, C.; Qin, J.; Dou, Q.; and Heng, P.-A. 2021. Feddg: Federated domain generalization on medical image segmentation via episodic learning in continuous frequency space. In *CVPR*, 1013–1023.
- Luo, Y.; Yin, L.; Bai, W.; and Mao, K. 2020. An appraisal of incremental learning methods. *Entropy*, 22(11): 1190.
- Mai, Z.; Li, R.; Jeong, J.; Quispe, D.; Kim, H.; and Sanner, S. 2022. Online continual learning in image classification: An empirical survey. *Neurocomputing*, 469: 28–51.
- Miao, Z.; Wang, Z.; Chen, W.; and Qiu, Q. 2022. Continual learning with filter atom swapping. In *ICLR*.
- Shin, H.; Lee, J. K.; Kim, J.; and Kim, J. 2017. Continual learning with deep generative replay. *NeurIPS*, 30.
- Su, Y.; Wang, X.; Qin, Y.; Chan, C.-M.; Lin, Y.; Wang, H.; et al. 2022. On transferability of prompt tuning for natural language processing. In *NAACL*, 3949–3969.
- Veličković, P.; Cucurull, G.; Casanova, A.; Romero, A.; Lio, P.; and Bengio, Y. 2017. Graph attention networks. *arXiv preprint arXiv:1710.10903*.
- Wang, Y.; Huang, Z.; and Hong, X. 2022. S-prompts learning with pre-trained transformers: An occam’s razor for domain incremental learning. *NeurIPS*, 35: 5682–5695.
- Wang, Z.; Zhang, Z.; Ebrahimi, S.; Sun, R.; Zhang, H.; Lee, C.-Y.; et al. 2022a. Dualprompt: Complementary prompting for rehearsal-free continual learning. In *ECCV*, 631–648. Springer.
- Wang, Z.; Zhang, Z.; Lee, C.-Y.; Zhang, H.; Sun, R.; Ren, X.; et al. 2022b. Learning to prompt for continual learning. In *CVPR*, 139–149.
- Zenke, F.; Poole, B.; and Ganguli, S. 2017. Continual learning through synaptic intelligence. In *ICML*, 3987–3995. PMLR.
- Zhao, Z.; Lin, H.; Chen, H.; and Heng, P.-A. 2019. PFA-ScanNet: Pyramidal feature aggregation with synergistic learning for breast cancer metastasis analysis. In *MICCAI*, 586–594. Springer.