

PM-DETR: Domain Adaptive Prompt Memory for Object Detection with Transformers

Peidong Jia*
 Jiaming Liu*
 Peking University
 China

Jiarui Wu
 Beihang University
 China

Senqiao Yang
 Harbin Institute of Technology, Shenzhen
 China

Xiaodong Xie
 Shanghang Zhang
 Peking University
 China

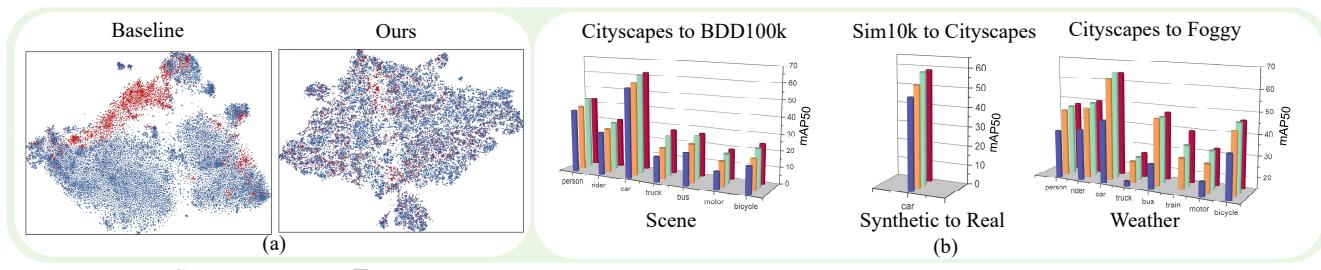


Figure 1: (a) compares the t-SNE results of different methods on the source and target domain data, and our method aligns the domain shift well compared to the baseline method. (b) indicates that our method achieves state-of-the-art (SOTA) performance on three challenging domain adaptation benchmarks.

ABSTRACT

The Transformer-based detectors (i.e., DETR) have demonstrated impressive performance on end-to-end object detection. However, transferring DETR to different data distributions may lead to a significant performance degradation. Existing adaptation techniques focus on model-based approaches, which aim to leverage feature alignment to narrow the distribution shift between different domains. In this study, we propose a hierarchical Prompt Domain Memory (PDM) for adapting detection transformers to different distributions. PDM comprehensively leverages the prompt memory to extract domain-specific knowledge and explicitly constructs a long-term memory space for the data distribution, which represents better domain diversity compared to existing methods. Specifically, each prompt and its corresponding distribution value are paired in the memory space, and we inject top M distribution-similar prompts into the input and multi-level embeddings of DETR. Additionally, we introduce the Prompt Memory Alignment (PMA) to reduce the discrepancy between the source and target domains by fully leveraging the domain-specific knowledge extracted from the prompt domain memory. Extensive experiments demonstrate that our method outperforms state-of-the-art domain adaptive object detection methods on three benchmarks, including scene, synthetic to real, and weather adaptation. Codes will be released.

1 INTRODUCTION

Object detection is a crucial computer vision task and serves as a prerequisite for various real-world applications, such as autonomous driving [1, 9, 25], visual grounding [35, 60], and manipulation [49, 64]. Convolutional Neural Networks (CNN) detectors [37, 42, 43] have shown satisfactory results, but they heavily depend on hand-crafted operations like non-maximum suppression. In recent times, a series of DEtection TRansformer (DETR) methods [7, 71] have been proposed with an end-to-end pipeline, which delivers promising performance when the test data is from the same distribution as the training data. However, such a fixed distribution is not typical in real-world scenarios [41], which often comprise diverse and disparate domains. When applying pre-trained DETR models, distribution shift commonly occurs [47], leading to significant performance degradation on the target data.

Existing adaptation techniques for DETR mainly rely on model-based approaches [53, 63], which aim to narrow the distribution shift between different domains via sequence feature alignment. Recent developments in prompt learning for both natural language processing (NLP) [31, 33, 36] and computer vision [26] have motivated researchers to introduce visual prompts in domain adaptation tasks. Several recent studies [8, 15, 19, 58] have leveraged prompts randomly set at the image or feature-level and fine-tuned them to extract domain-specific or maintain domain-invariant knowledge. These approaches offer a prompt-based perspective to address distribution shift, which can further aid the model-based methods in achieving better representations in the target domain. However, when prompt-based methods are applied in the target domain

*Both authors contributed equally to this research.

with various scene conversion and complex distribution data (i.e., autonomous driving data), the prompts are difficult to learn the long-term domain knowledge for the full data. Meanwhile, since object detection is an instance-level task and exists multiple objects in each sample, the previous prompt methods are hard to extract diverse domain knowledge for each category.

To this end, we propose a hierarchical Prompt Domain Memory (PDM) for adapting detection transformers, which can extract domain-specific knowledge by learning a set of prompts that dynamically instruct the DETR. Specifically, each prompt and its corresponding distribution value are paired in the memory space, and we dynamically select top M distribution similar prompts for each sample using the value. PDM explicitly constructs a long-term memory space for the detection transformer, allowing DETR to learn complex data distribution and different category domain knowledge at multi-levels, including input, token, and query levels. With the help of PDM, as shown in Fig.1 (a), feature representations in the two domains achieve smaller distribution shift compared to the previous method. However, while the prompt memory can extract more comprehensive domain-specific knowledge, it cannot reduce the distribution distance between different domains [15]. To address this limitation, we propose the Prompt Memory Alignment (PMA) method, which reduces the discrepancy between source and target domains in the Unsupervised Domain Adaptation (UDA) task. Traditional feature alignment methods [53, 63] can only align a small number of different domain samples in each iteration due to the limited GPU memory. Different from previous methods, since the proposed prompt memory can better represent the diversity of each domain, the PMA can fully leverage the domain-specific knowledge extracted from the memory and efficiently address the distribution shift. In addition, along with introducing PDM, we make the first attempt to design the visual prompt alignment strategy to jointly address the domain shift problem. In conclusion, our proposed approach of PDM and PMA enhances the performance of detection transformers in adapting to target domains by extracting diverse domain-specific knowledge and reducing the discrepancy between source and target domains.

We evaluate the prompt-based PM-DETR on three challenging benchmarks of UDA, including scene adaptation (Cityscapes [12] to BDD100k [62]), synthetic to real adaptation (Sim10k [27] to Cityscapes), and weather adaptation (Cityscapes to Foggy Cityscapes [46]). Our method outperforms state-of-the-art (SOTA) domain adaptive object detection methods, which improves the result to 58.6%, 33.3%, and 44.3% mAP in the three benchmarks, shown in Fig.1 (b).

The main contributions are summarized as follows:

- 1) We propose a hierarchical Prompt Domain Memory (PDM) to adapt detection transformers to different distributions, which constructs a long-term memory space to fully learn the complex data distribution and diversity domain-specific knowledge.
- 2) In order to better apply PDM in the Unsupervised Domain Adaptation (UDA), we propose the Prompt Memory Alignment (PMA) method to reduce the distribution distance between two domains, which can fully leverage the domain-specific knowledge extracted from the memory space.
- 3) We conduct extensive experiments on three challenging UDA scenarios to evaluate the effectiveness of our method. The method

achieves SOTA performance in all scenarios, including scene, synthetic to real, and weather adaptation.

2 RELATED WORKS

2.1 Object Detection

Object detection is a critical task of computer vision [21, 30, 65]. Previous convolutional neural network (CNN)-based approaches can be broadly categorized into two groups: the more complex two-stage methods [34, 43, 59] and the lighter one-stage methods [37, 42, 50]. However, these approaches exhibit a significant limitation due to their heavy reliance on handcrafted processes and initial guesses, particularly the non-maximum suppression (NMS) post-processing, which hinders their ability to be trained end-to-end. Recent advancements, such as DETR [7] and Deformable DETR [71], have addressed this issue by incorporating vision Transformers [52]. Deformable DETR introduces an innovative deformable multi-head attention mechanism that enables sparsity in attention and multi-scale feature aggregation without necessitating a feature pyramid structure. This innovation results in faster training and enhanced performance. For the critical issue of DETR, slow training convergence, Conditional DETR [39] speed up DETR training by leveraging a conditional cross-attention mechanism. In order to more effectively utilize the attention mechanism, SMCA [18] proposes a co-attention scheme that expedites DETR convergence, which includes multi-head and scale-selection attention. Moreover, DN-DETR [32] offers a novel perspective on faster training by employing a denoising approach to improve the stability of bipartite graph matching during the training stage. In our study, we use the classical Deformable DETR as the base detector and make the first attempt to introduce domain prompts into its workflow. We also design a hierarchical domain prompt memory to facilitate diverse domain knowledge extraction.

2.2 Domain adaptive object detection

Domain Adaptive Faster R-CNN [10] established the foundation for investigating domain-adaptive object detection techniques. Subsequent studies primarily utilized the adversarial training paradigm for cross-domain feature alignment. Various strategies have been proposed in these approaches to aggregate image or instance features, such as leveraging categorical predictions [56, 57] and exploiting spatial correlations [6, 57]. Hierarchical alignment of features was conducted at multiple levels, encompassing global, local, instance, and category levels [6, 38, 45, 56, 57]. Recent advancements in this field introduced innovative methods, including PICA [67], which specializes in few-shot domain adaptation, and Visually Similar Group Alignment (ViSGA)[44], which employs similarity-based hierarchical agglomerative clustering, achieving exceptional performance on specific benchmarks. Moreover, several studies explored alternative domain adaptation techniques or utilized different base detectors, such as Mean Teacher with Object Relations (MTOR)[6] and Unbiased Mean Teacher (UMT) [14]. Regarding the Transformer object detector, existing adaptation techniques for DETR predominantly rely on model-based approaches [53, 63], aiming to reduce the distribution shift between different domains through sequence feature alignment. In this paper, we provide a novel perspective on DETR cross-domain transfer by introducing a prompt-based approach. Specifically, we propose a Prompt-based Domain Memory

(PDM) and Prompt Pool Alignment (PMA) to enhance the performance of detection transformers in adapting to target domains. This is achieved by extracting diverse domain-specific knowledge and minimizing the discrepancy between source and target domains.

2.3 Prompt Learning

Prompt learning, originally introduced in the field of natural language processing (NLP), aims to adapt pre-trained language models to various downstream tasks in a parameter-efficient manner [5, 36, 40, 61, 68, 69]. Recently, researchers have extended the paradigm of prompt learning to efficient fine-tuning of vision models [2]. VPT [26] and its variants [11, 48] introduce minimal trainable parameters at the image or feature level of Transformer-based models for efficient transfer learning. L2P [55] and its follow-up method [54] propose a prompt pool-based approach for continual learning, aiming to avoid catastrophic forgetting and error accumulation. More recently, visual prompts have shown promising results in domain adaptation. DAPL [20] made the initial attempt to incorporate visual prompts into unsupervised domain adaptation (UDA). Subsequent studies, such as [8, 15, 19], explored diverse approaches to leverage visual prompts for classification domain adaptation problems. Additionally, SVDP [58] proposed a sparse visual prompt for efficient adaptation in segmentation tasks. However, these studies primarily focus on image-level [15] and pixel-level [58] domain adaptation tasks and are not optimized for instance-level transformer detection. Furthermore, when prompt-based methods are applied in the target domain with various scene conversions and complex data distributions [63] (e.g., autonomous driving data), it becomes challenging to utilize previous lightweight methods [15] to learn long-term domain knowledge. Due to the instance-level property, with multiple objects present in each sample, previous prompt methods [16] struggle to extract diverse domain knowledge for each category. To address this issue, we design a Prompt-based Domain Memory (PDM) tailored for adapting detection transformers, particularly in scenarios involving diverse data distributions.

3 METHODS

Preliminary This section introduces PM-DETR for transformer-based domain adaptive detectors. Given a model $D_S(y|x)$ trained in labeled source domain samples $\mathcal{D}_S(x, y)$, where x and y represent input data and Ground Truth, respectively. Our goal is adapting D_S to target model D_T through unlabeled target domain data $\mathcal{D}_T(x)$.

We propose a comprehensive prompt-based method to enhance the performance of detection transformers in adapting to target domains by extracting diverse domain-specific knowledge and reducing the discrepancy between source and target domains. In section 3.1, we first verify the motivation of the prompt-based adaptation method for transformer detectors. In section 3.2, a hierarchical Prompt Domain Memory (PDM) is illustrated to extract domain-specific knowledge in the multi-level transformer latent space. In section 3.3, Prompt Memory Alignment (PMA) is proposed to reduce the discrepancy between source and target domains. Finally, in Section 3.4, we elaborate on our training policy in detail. The overall pipeline is shown in Fig. 4, and Deformable DETR [71] is utilized as our default detection network.

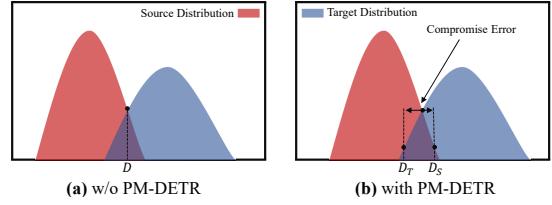


Figure 2: Demonstrate the compromise error. The red and blue areas represent the probability of correct classification in the source and target domains. When the models of the two domains are identical, model will converge to D . Our decoupling method converges to D_S and D_T respectively, avoiding the compromise error.

3.1 Motivation of Prompt-based Method

We verify the motivates of introducing a prompt-based method and constructing a long-term prompt domain memory. First, there is a brief explanation of how the traditional model-based adaptation method tackles the Unsupervised Domain Adaptation (UDA) problem. These methods [17, 24, 63] pursue to shrink upper boundary of the target error err_T by the sum of the source error err_S and a notion of distance d_H between the source and the target distributions in hypothesis space \mathcal{H} . Since err_S is primarily influenced by model complexity, researchers have predominantly focused on minimizing the inter-domain distance d_H . Suppose that construct a unified dataset \mathbb{U} defined as below:

$$\mathbb{U} = \{x_i, y=0\}_{i=1}^p \cup \{x_j, y=1\}_{j=p+1}^q \quad (1)$$

Where the first p samples are from source domain \mathcal{D}_S and labeled as 0, the rest samples are from target domain \mathcal{D}_T and labeled as 1. By constructing this unified dataset \mathbb{U} , we create a share high dimensional space for both the source and target domains, allowing us to effectively measure the distance between their distributions in the hypothesis space \mathcal{H} , shown in Fig. 2 (a). This unified dataset serves as a foundation for minimizing the inter-domain distance and enabling better adaptation from the source to the target domain. Furthermore, the work of Ben-David et al. [3, 4] has provided evidence that the empirical \mathcal{H} -divergence between two domains can be computed by the following equation:

$$\begin{aligned} d_H(S, T) &= 2(1 - \min_{D \in \mathcal{H}} [\frac{1}{p} \sum_{i=1}^p D[\mathbb{U}(y_i = 0)] \\ &\quad + \frac{1}{q-p} \sum_{j=p+1}^{q-p} D[\mathbb{U}(y_j = 1)]]]) \\ &= 2(1 - \varepsilon_D^S - \varepsilon_D^T) \end{aligned} \quad (2)$$

The methods discussed above are based on an intuitive assumption that model parameters capable of fitting both the source and target domain data well can be learned in the same hypothesis space \mathcal{H} . However, it is important to note that this assumption does not always hold true in practical scenarios. In reality, there can be inherent differences between the source and target domains that make it challenging to find a single hypothesis space that adequately captures both domains. As a result, when attempting to adapt a model from the source to the target domain, a compromise error may be introduced. Fig. 2 visually illustrates this compromise

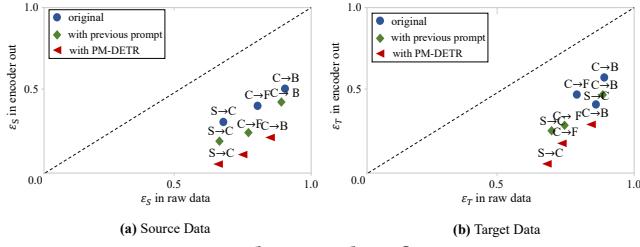


Figure 3: Quantitative domain classification errors comparison of baseline and our method. Where $C \rightarrow B$ means from Cityscapes to BDD100k, $S \rightarrow C$ means from Sim10k to Cityscapes, $C \rightarrow F$ means from Cityscapes to Foggy.

error, which represents the discrepancy between the optimal models for each domain and the compromise model that attempts to accommodate both domains.

To this end, drawing inspiration from soft prompt learning techniques used in NLP and computer vision tasks, where the pre-trained model is adapted to different downstream tasks through prompt manipulation within the input sequence, we propose that incorporating a lightweight visual prompt can assist in bounding the inter-domain distance $d_{\mathcal{H}}$ within a smaller interval. Specifically, we adopt domain prompt warp into input image, encoder embedding, and decoder queries, so that the hypothesis space in source and target domain can be decoupled to \mathcal{H}_S and \mathcal{H}_T . The prompt memory explicitly constructs a long-term memory space for better representing the diversity of domain knowledge, which further assist the hypothesis space decoupling. Under different hypothesis spaces, prompt memory alignment encourages mining in-domain knowledge by constraints, so the model is optimized as D_S and D_T in the source and target domains, respectively, as shown in Fig. 2 (b). Thus, in theory, the inter-domain distance will be reduced by the following equation

$$d_{\mathcal{H}_S, \mathcal{H}_T}(S, T) = 2(1 - \epsilon_{D_S}^S - \epsilon_{D_T}^T) < d_{\mathcal{H}}(S, T) \quad (3)$$

To further substantiate the presence of compromise error, we conduct experiments on three domain adaptive object detection tasks as depicted in Fig. 3. The figure quantitatively compares the classification errors of the model-based, previous prompt-based, and Prompt Domain Memory (PDM) method on a domain classifier consisting of three multi-layer perceptrons in series. As we can see, the previous prompt-based method achieve smaller classification errors compared to the model-based method, and the proposed PDM further has a significant improvement in classification errors. As proven by [17], generalization upper boundary on the target risk can be smaller attributed to lower classification error, which means that model will perform better in the target domain.

3.2 Hierarchical Prompt Domain Memory

We propose prompt domain memory P_S and P_T for source and target domains, respectively. Each memory pool P has N prompt pairs $\{< v_i, p_i > | i = 1, \dots, N, v \in \mathbb{R}^{1 \times d}, p \in \mathbb{R}^{L \times d}\}$, where v indicates prompt distribution value and p indicates visual prompt weight. L stands for embedding length and d stands for embedding dimension. We randomly initialize v and p . Visual prompt weight p stands for the aggregation of domain-specific knowledge, which

can be learned by warping in the prefix of the input sequence. Prompt distribution value v will be used to measure the distribution similarity with the input. In this way, prompt memory can cover diverse domain knowledge for complex data distribution.

Hierarchical Warp Position. We introduce prompt memory in three crucial embeddings including input image, encoder token, and decoder queries. The motivations are three-fold. First, the combination of multi-level prompts provides a comprehensive domain transfer mechanism, where different levels of prompts decrease domain shift that could not be narrowed by previous levels. Second, prompt in decoder query can align the distribution of objects in different domain datasets, thus improving the recall of Deformable DETR. Third, the multi-level prompts only increase the number of parameters by a very small amount (0.063% of model parameters), but it can greatly enhance the plasticity of the model and release the power in learning diverse domain-specific representations. Our ablation experiments in Sec. 4.2 show that the hierarchical prompt memory can better represent the domain diversity and boost the performance of Deformable DETR in the target domain.

Distribution Value Similarity Selection. The image scenes, object classes, and object distributions in the target dataset have large variances, which often leads to sub-optimal performance if only use the same prompt for each instance to extract domain knowledge. We design a distribution-guided strategy to adaptively select prompts from prompt memory. We project input embedding by transformation function γ to v 's shape. Here we utilize the average mean along the embedding channel to aggregate input characteristics. Then we calculate the cosine similarity between v and the projection embedding utilizing function ψ .

$$V_M = \operatorname{argmax} \sum_{i=1}^M \psi(V, \gamma(x)) \quad (4)$$

According to the cosine similarity value, V_M is the selected nearest M neighbor prompts (i.e. $M = 4$) in prompt memory.

3.3 Prompt Memory Alignment

After receiving the domain knowledge transferred from prompt memory, we further introduce Prompt Memory Alignment (PMA) to address the domain shift accumulation. For the encoder phase, we aim to pull close the input and token level prompts from two domain prompt memory, as shown in Fig. 4. Specifically, we utilize respective MLPs to project prompt tokens to a shared embedding space, in which the dimension is $L \times C \times 2$, L equals to encoder token length, channel dimension C is set to 256. We adopt encoder prompt alignment loss \mathcal{L}_{epa} to pull close the two domain prompt embeddings and explicitly constrain prompt to learn in-domain knowledge, as shown in Eq. 5.

$$\mathcal{L}_{epa}(X, D) = \lambda_1 \min D(X_{i<|p|\times M}) + \lambda_2 \max D(X_{i\geq|p|\times M}) \quad (5)$$

Where D denotes the domain discriminator. For decoder phase, decoder prompt alignment loss \mathcal{L}_{dpa} , similar to \mathcal{L}_{epa} , are proposed. Since objects in the same category and spatially connected tend to be visually similar, \mathcal{L}_{dpa} constraint prompts in decoder queries to learn from categorical and spatial correlations while decreasing data distribution distance between the two domains.

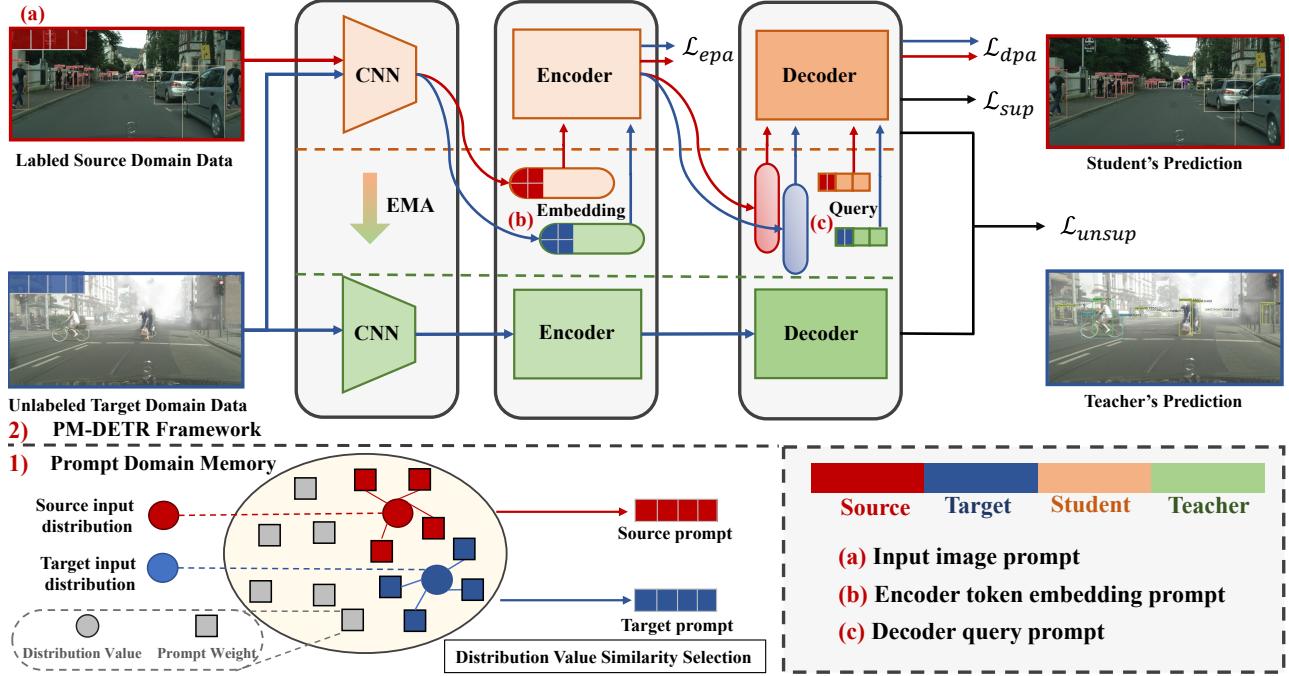


Figure 4: The Overall framework of PM-DETR. (1) **Prompt Domain Memory (PDM).** Hierarchical prompt domain memory, which warps on input, encoder token, and decoder query, learns a set of prompts for excavating diverse domain-specific knowledge. We adaptively select the top M prompts by distribution value similarity. The selection strategy as described in Eq. 4. (2) **PM-DETR Framework.** We construct a teacher-student paradigm to optimize the student model and PDM by source domain labels and target domain pseudo labels. Besides, we propose a novel prompt memory alignment (PMA) to constraint multi-level prompts digging in-domain knowledge and spatial correlations, as described in Eq. 5.

3.4 Overall Optimization for PM-DETR

PM-DETR leverages the teacher-student framework, which includes two models with the same architecture and weights at initialization. During training, the student model is updated using back-propagation, while the teacher model is updated by taking the Exponential Moving Average (EMA) of the student’s weights. The weights of the teacher model θ'_t at time step t is calculated by taking a weighted average of the teacher’s previous weights and the current student’s weights θ_t :

$$\theta'_t = \alpha\theta'_{t-1} + (1 - \alpha)\theta_t \quad (6)$$

α is a smoothing coefficient hyperparameter and is set to 0.999. And the object query embeddings are kept the same between the two models to enhance consistency. The object queries are trainable embeddings, initialized with the normal distribution at the start of the training procedure. Then we use the temporally ensembled teacher model to optimize the parameter of the student model and prompts in the target domain via pseudo labels.

For the integral optimizing, the first loss is a penalty on the source domain for supervised learning (\mathcal{L}_{sup}), which distills domain independent generic features and avoids catastrophic forgetting. The second loss uses the target domain pseudo-labels generated by the teacher model for unsupervised learning (\mathcal{L}_{unsup}) to extract the target domain knowledge. The two losses are separately used to optimize source and target domain prompt memory. Combined

with the proposed prompt alignment losses, the overall constraint function is:

$$\mathcal{L} = \lambda_s \mathcal{L}_{sup} + \lambda_{us} \mathcal{L}_{unsup} + \lambda_{epa} \mathcal{L}_{epa} + \lambda_{dpa} \mathcal{L}_{dpa} \quad (7)$$

To maintain the balance of loss penalties, λ_s and λ_{us} are set to 1, λ_{epa} , and λ_{dpa} are set to 0.25. The detection loss (\mathcal{L}_{sup} and \mathcal{L}_{unsup}) are combined by focal loss and L1 loss [7, 71].

4 EXPERIMENTS

We conduct extensive experiments to demonstrate the advantages of our proposed method for object detection Unsupervised Domain Adaptation (UDA) task. In Section 4.1, we provide a description of the datasets, as well as the details of model settings. In Section 4.2, we compare mean Average Precision (mAP) metric of PM-DETR with other baselines [6, 23, 28, 43–45, 50, 53, 57, 63, 70, 71] in three challenging domain adaptation scenarios, including Weather, Scene, and Synthetic to Real Adaptation. Comprehensive ablation studies are conducted to investigate the impact of each component in Section 4.3. Furthermore in Section 4.4, qualitative analysis is given to facilitate intuitive understanding.

4.1 Experimental Setup

Datasets. We evaluate our method on four public datasets, including Cityscapes [12], Foggy Cityscapes [46], Sim10k [66], and BDD100k [62]. These datasets provide diverse and challenging scenarios for domain adaptation tasks:

Table 1: Performance comparison of different methods for weather adaptation, that is, from Cityscapes to Foggy Cityscapes. FRCNN and DefDETR are abbreviations for Faster R-CNN and Deformable DETR, respectively.

Method	Detector	Publication	person	rider	car	truck	bus	train	mcycle	bicycle	mAP	Gain
<i>Two Stage :</i>												
FasterRCNN [43](Source)	FRCNN	NIPS2015	26.9	38.2	35.6	18.3	32.4	9.6	25.8	28.6	26.9	00.00
CR-DA [56]	FRCNN	CVPR2020	30.0	41.2	46.1	22.5	43.2	27.9	27.8	34.7	34.2	+07.3
DivMatch [28]	FRCNN	CVPR2019	31.8	40.5	51.0	20.9	41.8	34.3	26.6	32.4	34.9	+08.0
MTOR [6]	FRCNN	CVPR2019	30.6	41.4	44.0	21.9	38.6	40.6	28.3	35.6	35.1	+08.2
SWDA [45]	FRCNN	CVPR2019	31.8	44.3	48.9	21.0	43.8	28	28.9	35.8	35.3	+08.4
SCDA [70]	FRCNN	CVPR2019	33.8	42.1	52.1	26.8	42.5	26.5	29.2	34.5	35.9	+09.0
CR-SW [56]	FRCNN	CVPR2020	34.1	44.3	53.5	24.4	44.8	38.1	26.8	34.9	37.6	+10.7
GPA [57]	FRCNN	CVPR2020	32.9	46.7	54.1	24.7	45.7	41.1	32.4	38.7	39.5	+12.6
ViSGA [44]	FRCNN	ICCV2021	38.8	45.9	57.2	29.9	50.2	51.9	31.9	40.9	43.3	+16.4
<i>One Stage :</i>												
FCOS [50] (Source)	FCOS	ICCV2019	36.9	36.3	44.1	18.6	29.3	8.4	20.3	31.9	28.2	+01.3
EPM[23]	FCOS	ECCV2020	44.2	46.6	58.5	24.8	45.2	29.1	28.6	34.6	39.0	+12.1
<i>Transformer based :</i>												
Def DETR [71] (Source)	DefDETR	ICLR2021	37.7	39.1	44.2	17.2	26.8	5.8	21.6	35.5	28.5	+01.6
SFA [53]	DefDETR	ACMMM2021	46.5	48.6	62.6	25.1	46.2	29.4	28.3	44.0	41.3	+14.4
MTTrans[63]	DefDETR	ECCV2022	47.7	49.9	65.2	25.8	45.9	33.8	32.6	46.5	43.4	+16.5
Ours(PM-DETR)	DefDETR	-	47.8	50.2	64.7	26.5	47.2	39.6	32.4	46.1	44.3	+17.4

Table 2: Performance comparison of different methods for the scene adaptation, i.e., Cityscapes to BDD100k daytime subset.

Methods	Detector	Publication	person	rider	car	truck	bus	mcycle	bicycle	mAP	Gain
<i>Two Stage :</i>											
FasterRCNN [43](Source)	FRCNN	NIPS2015	28.8	25.4	44.1	17.9	16.1	13.9	22.4	24.1	0.00
DAF [10]	FRCNN	CVPR2018	28.9	27.4	44.2	19.1	18.0	14.2	22.4	24.9	+0.8
SCDA [70]	FRCNN	CVPR2019	29.3	29.2	44.4	20.3	19.6	14.8	23.2	25.8	+1.7
CR-DA [56]	FRCNN	CVPR2020	30.8	29.0	44.8	20.5	19.8	14.1	22.8	26.0	+1.9
SWDA [45]	FRCNN	CVPR2019	29.5	29.9	44.8	20.2	20.7	15.2	23.1	26.2	+2.1
CR-SW [56]	FRCNN	CVPR2020	32.8	29.3	45.8	22.7	20.6	14.9	25.5	27.4	+3.3
<i>One Stage :</i>											
FCOS [50](Source)	FCOS	ICCV2019	38.6	24.8	54.5	17.2	16.3	15.0	18.3	26.4	+2.3
EPM [23]	FCOS	ECCV2020	39.6	26.8	55.8	18.8	19.1	14.5	20.1	27.8	+3.7
<i>Transformer Based :</i>											
Def DETR [71](Source)	DefDETR	ICLR2021	38.9	26.7	55.2	15.7	19.7	10.8	16.2	26.2	+2.1
SFA [53]	DefDETR	ACMMM2021	40.2	27.6	57.5	19.1	23.4	15.4	19.2	28.9	+4.8
MTTrans [63]	DefDETR	ECCV2022	44.1	30.1	61.5	25.1	26.9	17.7	23.0	32.6	+8.5
Ours(PM-DETR)	DefDETR	-	43.3	30.9	62.0	27.4	26.8	18.7	23.9	33.3	+9.2

Weather Adaptation. In this scenario, we use Cityscapes as the source dataset, consisting of 2,975 training images and 500 evaluation images. The target dataset is Foggy Cityscapes, generated from Cityscapes using a fog synthesis algorithm. Foggy Cityscapes introduces foggy conditions to the images, enabling us to evaluate the performance of our method in adapting object detection models from clear weather to foggy weather scenarios.

Scene Adaptation. In this condition, Cityscapes serves as the source dataset, while the target dataset is the daytime subset of BDD100k. BDD100k consists of 36,728 training images and 5,258 validation images, all annotated with bounding boxes. This subset provides a diverse range of scenes captured during the daytime

Synthetic to Real Adaptation. In this particular scenario, we employ Sim10k as the source domain, which is generated using the Grand Theft Auto game engine. Sim10k comprises 10,000 training images, accompanied by 58,701 bounding box annotations. As for the target domain, we utilize the car instances from Cityscapes for both training and evaluation purposes.

Implementation Details. Our method is built based on Deformable DETR [71]. We set ImageNet [13] pre-trained ResNet-50 [22] as CNN backbone in all experiments. In the burn-in step, we adopt Adam optimizer [29] for training over 50 epochs. The initial learning rate is set to $2e - 04$, which decayed by 0.1 after 40 epochs. The batch size is set to 4 for all adaptation scenarios. In the second cross-domain training step, the model is trained for 12 epochs. The prompt-based parameters are initialized with random float numbers and set the learning rate to $2e - 05$. The learning rate for the other parameters, excluding the prompt-based parameters, is relatively

small and set to $2e - 06$. All learning rates decayed by 0.1 after 10 epochs. In addition, we adopt mean Average Precision (mAP) with a threshold of 0.5 as the evaluation metric. We set the filtering threshold (confidence) for the pseudo-label generation to 0.5. All experiments are conducted on two NVIDIA Tesla A100 GPUs.

4.2 Comparisons with SOTA Methods

Weather Adaptation. To assess the reliability of object detectors under varying weather conditions, we transfer models from Cityscapes to Foggy Cityscapes. As shown in Table 1, our proposed method PM-DETR significantly outperforms other cutting-edge approaches, achieving a 44.3% score compared to the closest SOTA end-to-end model, MTTrans [63], at 43.4%. Furthermore, it reveals that PM-DETR considerably enhances Deformable DETR's cross-domain performance, achieving a 15.8% absolute gain in mAP50 and outperforming all previous domain adaptive object detection methods. These promising results highlight the ability of our method to extract diverse domain-specific knowledge and effectively address distribution shift, leading to improved performance in unsupervised domain adaptation for object detection tasks.

Scene Adaptation. In real-world applications, such as autonomous driving, scene layouts are not static and frequently change. It makes model performance under scene adaptation crucial. Our proposed method, PM-DETR, demonstrates its effectiveness in scene adaptation as shown in Table 2, achieving SOTA results (33.3%) and significantly improving upon previous works. Additionally, the performance of five out of seven categories in the target domain dataset has been enhanced.

Table 3: Performance comparison of different methods for the synthetic to real adaptation, i.e. Sim10k to Cityscapes.

Methods	Detector	Publication	mAP(car)	Gain
FasterRCNN [43](Source)	FRCNN	NIPS2015	34.6	00.00
DAF [10]	FRCNN	CVPR2018	41.9	+07.3
CR-DA [56]	FRCNN	CVPR2020	43.1	+08.5
DivMatch [28]	FRCNN	CVPR2019	43.9	+09.3
SWDA [45]	FRCNN	CVPR2019	44.6	+10.0
SCDA [70]	FRCNN	CVPR2019	45.1	+10.5
CR-SW [56]	FRCNN	CVPR2020	46.2	+11.6
MTOR [6]	FRCNN	CVPR2019	46.6	+12.0
GPA [57]	FRCNN	CVPR2020	47.6	+13.0
ViSGA [44]	FRCNN	ICCV2021	49.3	+14.7
FCOS [50](Source)	FCOS	ICCV2019	42.5	+7.9
EPM [23]	FCOS	ECCV2020	47.3	+12.7
DefDETR [71](Source)	DefDETR	ICLR2021	47.4	+12.8
SFA [53]	DefDETR	ACMMM2021	52.6	+23.3
MTTrans [63]	DefDETR	ECCV2022	57.9	+23.3
PM-DETR(Ours)	DefDETR	-	58.6	+24.0

Synthetic to Real Adaptation. The training process of object detectors using affordable and accurate simulation datasets has been proven to yield improved performance. However, this approach also brings about a notable challenge in the form of a significant inter-domain gap. In the synthetic to real adaptation scenario, we evaluated the performance of our proposed method, PM-DETR, as shown in Table 3. PM-DETR achieved state-of-the-art accuracy with a mAP of 58.6%, outperforming Deformable DETR by 11.2% mAP. These promising results further demonstrate the importance of a long-term domain memory space for transformer detectors to effectively extract comprehensive domain knowledge in real-world unsupervised domain adaptation scenarios.

4.3 Ablation Study

Effectiveness of each component. To better analyze each component in our proposed PM-DETR framework, we conduct ablation studies by accruing parts of the components in PM-DETR. As presented in Table 4 (PM-DETR-AS0), the teacher-student structure is a common technique in UDA [6, 63], which is used to generate pseudo labels in the target domain and has 8.5% mAP drop compared to our method. This verifies the improvement of our method does not come from the usage of this prevalent scheme and the model still suffers from the domain shift problem due to imperfect target domain feature extraction. In PM-DETR-AS11, by introducing prompt domain memory (PDM) in the input image, we observe that the mAP increase by 7.3%. When employing PDM in encoder token embedding (PM-DETR-AS12) and in decoder query embedding (PM-DETR-AS13), mAP improves by 6.5% and 6.9%, respectively. The result clearly demonstrates that the utilization of a long-term memory space enables the model to fully learn the complex data distribution and capture diverse domain-specific knowledge in multiple levels of DETR. In terms of Prompt Memory Alignment (PMA), PM-DETR-AS21 improves the mAP to 43.8% by encoder prompt alignment, and PM-DETR-AS22 improves the mAP 43.9% by decoder prompt alignment. The improved performance evaluates that PMA can further reduce the discrepancy between the two domains. PM-DETR shows the complete combination of all components which achieves 15.8% improvement in total. It proves that all components compensate each other and jointly mitigate the object detection domain shift problem in an unsupervised paradigm.

How do the prompt memory size and selection strategy affect the performance? In Fig. 6 (a), we observe the impact of

Table 4: Ablation studies on the weather adaptation scenario. MT stands for the mean teacher framework. Img. and Emd. stand for the input image and encoder embedding. PDM and PMA are the abbreviations of Prompt Domain Memory and Prompt Memory Alignment. Components of other experiments that differ from PM-DETR are marked in red.

Methods	MT	PDM			PMA		mAP50
		Img.	Emd.	Query	\mathcal{L}_{epa}	\mathcal{L}_{dpa}	
Def. DETR (Source)		✗	✗	✗	✗	✗	28.500
PM DETR-AS0	✓	✗	✗	✗	✗	✗	35.843
PM DETR-AS11	✓	✓	✗	✗	✗	✗	43.131
PM DETR-AS12	✓	✗	✓	✗	✗	✗	42.417
PM DETR-AS13	✓	✗	✗	✓	✗	✗	42.765
PM DETR-AS21	✓	✓	✓	✓	✓	✗	43.812
PM DETR-AS22	✓	✓	✓	✓	✗	✓	43.943
PM DETR	✓	✓	✓	✓	✓	✓	44.288

different prompt memory sizes on model performance. When using a single prompt (memory size equals one), there is a significant drop in performance, suggesting that a single prompt suffers severe diversity interference. As the memory size increases, the performance of the prompt-based model gradually improves and reaches its peak when the memory size is 10. Further increasing the size leads to a slight decrease in performance, but it still outperforms the single prompt scenario. This demonstrates that our prompt domain memory effectively constructs a long-term domain memory, enabling the model to understand complex data distribution and diverse domain-specific knowledge. Fig. 6 (b) compare the effect of different prompt selection schemes on model performance, including random, k-means, and distribution-based approaches. Our distribution-based selection strategy achieves the highest performance, highlighting the effectiveness of our method in capturing crucial domain-specific information. A selection method that considers the instance-level inputs can effectively stably handle the variance of the data distribution in the target domain. To better understand the prompt selection mechanism, we plot the prompt selection frequency histograms for three domain adaptive tasks in Fig. 5 (b). Our prompt selection mechanism clearly encourages more knowledge sharing between similar categories and more knowledge comparison between dissimilar categories.

4.4 Visualization and Analysis

Detection Results. We show some visualization results of PM-DETR on three target domain datasets, i.e. Foggy Cityscapes, BDD100k, and Cityscapes, accompanied by ground truth and previous state-of-the-art (SOTA) methods. As shown in Fig. 5 (a) Row1 (Cityscapes to Foggy Cityscapes), PM-DETR has higher recall and more accurate classification results in dense fog occlusion. As shown in Fig. 5 (a) Row2 (Cityscapes to BDD100k), our method properly classifies and locates objects even when they are heavily occluded or challengingly small in size. In Fig. 5 (a) Row3 (Sim10k to Cityscapes), we can even alleviate label misalignment (car & truck) without supervision to some degree. All visual results are consistent with the numerical assessment results in three target domains, indicating that PM-DETR manages to mitigate the domain shift problem in the UDA Transformer detector.

t-SNE Distribution Results. Following the t-distributed stochastic neighbor embedding (t-SNE) method [51], in Fig. 1 (a) and Fig. 5 (c), we visualize two types of t-SNE plots to illustrate the effectiveness

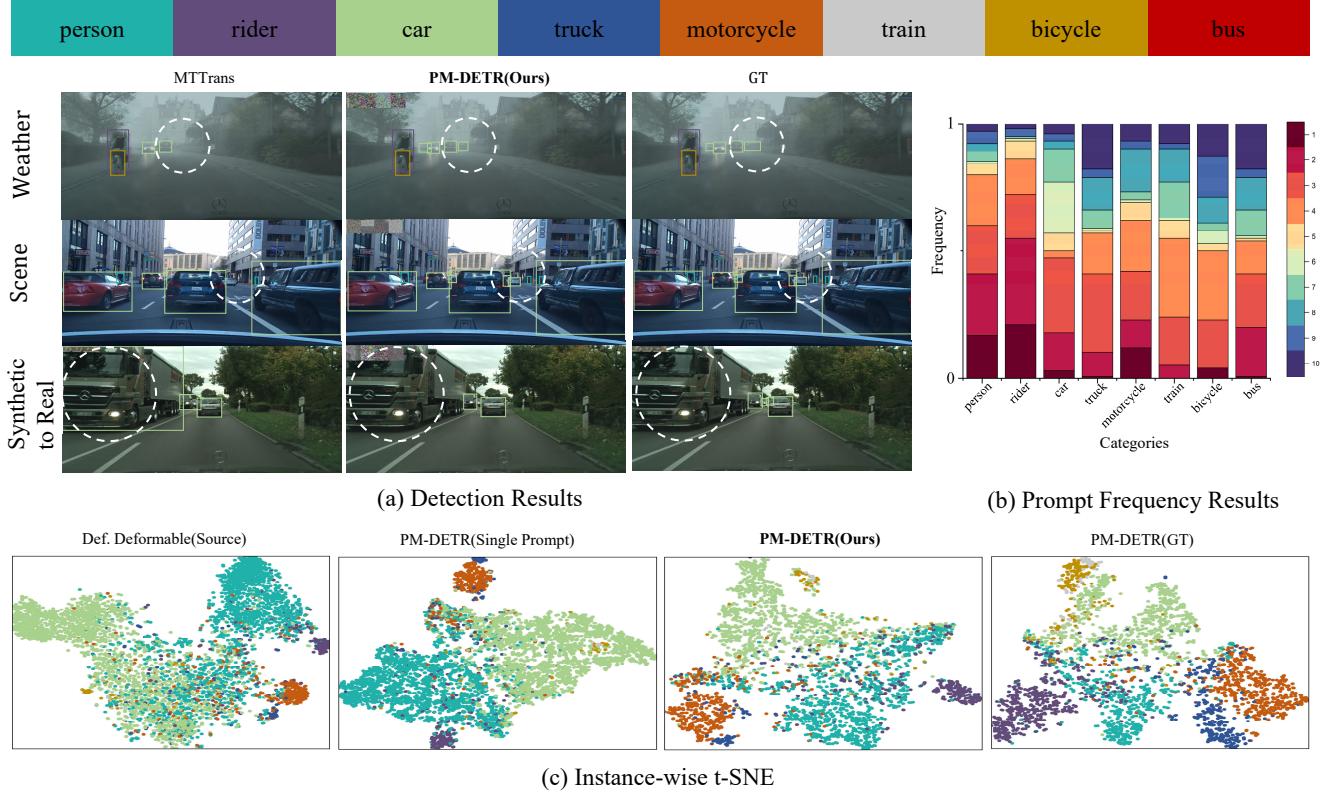


Figure 5: (a) Qualitative comparison of PM-DETR with previous SOTA method and GT in three scenarios. The white circle area reflects the superiority of our method. (b) In Cityscapes to Foggy Cityscapes scenario, the frequency statistics of prompt picking in memory space correspond to different categories. (c) In Cityscapes to Foggy Cityscapes scenario, instance-level feature t-SNE results. GT stands for training in Foggy Cityscapes within fully supervised learning.

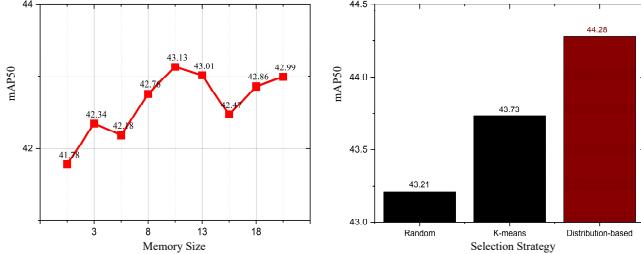


Figure 6: (a) Effects of prompts' memory size and (b) Selection strategy on the UDA task.

of our approach: global-wise t-SNE and instance-wise t-SNE. By examining the global-wise t-SNE plot, we can gain insights into how well our method mixes different domains. On the other hand, in the instance-wise t-SNE, we focus on visualizing individual instances and their embeddings. It can be observed that the t-SNE of our method is most similar to the results of fully supervised training in terms of inter-category distance as well as similar category aggregation. This firmly corroborates the ability of our method in mining diverse domain-specific knowledge for each category.

5 CONCLUSIONS

This paper presents a novel prompt-based method to enhance the adaptation ability of transformer detection by decoupling hypothesis space and mitigating the existing compromise error. Our approach leverages a hierarchical Prompt Domain Memory (PDM) to maintain a long-term memory space that facilitates comprehensive learning of the complex data distribution and diverse domain-specific knowledge. To effectively utilize PDM in cross-domain learning, we propose the Prompt Memory Alignment (PMA) method, which reduces the distribution distance between two domains by boundedly extracting the domain-specific knowledge from the memory space. We evaluate the effectiveness of our method through extensive experiments on three challenging Unsupervised Domain Adaptation (UDA) scenarios. The results demonstrate the significant improvements achieved by our approaches, PDM and PMA jointly address the distribution shift problem. Moreover, our method is applicable across different domain distances, making it a versatile solution for various domain adaptation problems.

REFERENCES

- [1] Eduardo Arnold, Omar Y Al-Jarrah, Mehrdad Dianati, Saber Fallah, David Oxtoby, and Alex Mouzakitis. 2019. A survey on 3d object detection methods for autonomous driving applications. *IEEE Transactions on Intelligent Transportation Systems* 20, 10 (2019), 3782–3795.
- [2] Hyojin Bahng, Ali Jahanian, Swami Sankaranarayanan, and Phillip Isola. 2022. Exploring Visual Prompts for Adapting Large-Scale Models.
- [3] Shai Ben-David, John Blitzer, Koby Crammer, Alex Kulesza, Fernando Pereira, and Jennifer Wortman Vaughan. 2010. A theory of learning from different domains. *Machine learning* 79, 1 (2010), 151–175.
- [4] Shai Ben-David, John Blitzer, Koby Crammer, and Fernando Pereira. 2006. Analysis of representations for domain adaptation. *Advances in neural information processing systems* 19 (2006).
- [5] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language Models are Few-Shot Learners.
- [6] Qi Cai, Yingwei Pan, Chong-Wah Ngo, Xinmei Tian, Lingyu Duan, and Ting Yao. 2019. Exploring object relation in mean teacher for cross-domain detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 11457–11466.
- [7] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. 2020. End-to-end object detection with transformers. In *European conference on computer vision*. Springer, 213–229.
- [8] Haoran Chen, Zuxuan Wu, and Yu-Gang Jiang. 2022. Multi-Prompt Alignment for Multi-source Unsupervised Domain Adaptation. *arXiv preprint arXiv:2209.15210* (2022).
- [9] Xiaozhi Chen, Huimin Ma, Ji Wan, Bo Li, and Tian Xia. 2017. Multi-view 3d object detection network for autonomous driving. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*. 1907–1915.
- [10] Yuhua Chen, Wen Li, Christos Sakaridis, Dengxin Dai, and Luc Van Gool. 2018. Domain adaptive faster r-cnn for object detection in the wild. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 3339–3348.
- [11] Jonathan Conder, Josephine Jefferson, Nathan Pages, Khurram Jawed, Alireza Nejati, and Mark Sagar. 2022. Efficient Transfer Learning for Visual Tasks via Continuous Optimization of Prompts.
- [12] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. 2016. The cityscapes dataset for semantic urban scene understanding. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 3213–3223.
- [13] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. 2009. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*. Ieee, 248–255.
- [14] Jinzhong Deng, Wen Li, Yuhua Chen, and Lixin Duan. 2021. Unbiased mean teacher for cross-domain object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 4091–4101.
- [15] Yulu Gan, Xiancheng Ma, Yihang Lou, Yan Bai, Renrui Zhang, Nian Shi, and Lin Luo. 2022. Decorate the Newcomers: Visual Domain Prompt for Continual Test Time Adaptation. *arXiv preprint arXiv:2212.04145* (2022).
- [16] Yulu Gan, Mingjie Pan, Rongyu Zhang, Zijian Ling, Lingran Zhao, Jiaming Liu, and Shanghang Zhang. 2022. Cloud-Device Collaborative Adaptation to Continual Changing Environments in the Real-world. *arXiv preprint arXiv:2212.00972* (2022).
- [17] Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François Laviolette, Mario Marchand, and Victor Lempitsky. 2016. Domain-adversarial training of neural networks. *The journal of machine learning research* 17, 1 (2016), 2096–2030.
- [18] Peng Gao, Minghang Zheng, Xiaogang Wang, Jifeng Dai, and Hongsheng Li. 2021. Fast convergence of detr with spatially modulated co-attention. In *Proceedings of the IEEE/CVF international conference on computer vision*. 3621–3630.
- [19] Yunhe Gao, Xingjian Shi, Yi Zhu, Hao Wang, Zhiqiang Tang, Xiong Zhou, Mu Li, and Dimitris N Metaxas. 2022. Visual Prompt Tuning for Test-time Domain Adaptation. *arXiv preprint arXiv:2210.04831* (2022).
- [20] Chunjiang Ge, Rui Huang, Mixue Xie, Zihang Lai, Shiji Song, Shuang Li, and Gao Huang. 2022. Domain Adaptation via Prompt Learning. *ArXiv abs/2202.06687* (2022).
- [21] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. 2017. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*. 2961–2969.
- [22] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 770–778.
- [23] Cheng-Chun Hsu, Yi-Hsuan Tsai, Yen-Yu Lin, and Ming-Hsuan Yang. 2020. Every pixel matters: Center-aware feature alignment for domain adaptive object detector. In *European Conference on Computer Vision*. Springer, 733–748.
- [24] Shishuai Hu, Zehui Liao, and Yong Xia. 2022. ProSFDA: Prompt Learning based Source-free Domain Adaptation for Medical Image Segmentation. *arXiv preprint arXiv:2211.11514* (2022).
- [25] Yu Huang and Yue Chen. 2020. Autonomous driving with deep learning: A survey of state-of-art technologies. *arXiv preprint arXiv:2006.06091* (2020).
- [26] Menglin Jia, Luming Tang, Bor-Chun Chen, Claire Cardie, Serge Belongie, Bharath Hariharan, and Ser-Nam Lim. 2022. Visual prompt tuning. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXXIII*. Springer, 709–727.
- [27] Matthew Johnson-Roberson, Charles Barto, Rounak Mehta, Sharath Nittur Sridhar, Karl Rosaen, and Ram Vasudevan. 2016. Driving in the matrix: Can virtual worlds replace human-generated annotations for real world tasks? *arXiv preprint arXiv:1610.01983* (2016).
- [28] Taekyung Kim, Minki Jeong, Seunghyeon Kim, Seocheon Choi, and Changick Kim. 2019. Diversify and match: A domain adaptive representation learning paradigm for object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 12456–12465.
- [29] Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014).
- [30] Alexander Kirillov, Kaiming He, Ross Girshick, Carsten Rother, and Piotr Dollár. 2019. Panoptic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 9404–9413.
- [31] Brian Lester, Rami Al-Rfou, and Noah Constant. 2021. The power of scale for parameter-efficient prompt tuning. *arXiv preprint arXiv:2104.08691* (2021).
- [32] Feng Li, Hao Zhang, Shilong Liu, Jian Guo, Lionel M Ni, and Lei Zhang. 2022. Dn-detr: Accelerate detr training by introducing query denoising. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 13619–13627.
- [33] Xiang Lisa Li and Percy Liang. 2021. Prefix-tuning: Optimizing continuous prompts for generation. *arXiv preprint arXiv:2101.00190* (2021).
- [34] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. 2017. Feature pyramid networks for object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2117–2125.
- [35] Daqing Liu, Hanwang Zhang, Feng Wu, and Zheng-Jun Zha. 2019. Learning to assemble neural module tree networks for visual grounding. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 4673–4682.
- [36] Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. 2023. Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing. *Comput. Surveys* 55, 9 (2023), 1–35.
- [37] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott Reed, Cheng-Yang Fu, and Alexander C Berg. 2016. Ssd: Single shot multibox detector. In *European conference on computer vision*. Springer, 21–37.
- [38] Zhipeng Luo, Zhongang Cai, Changqing Zhou, Gongjie Zhang, Haiyu Zhao, Shuai Yi, Shijian Lu, Hongsheng Li, Shanghang Zhang, and Ziwei Liu. 2021. Unsupervised domain adaptive 3d detection with multi-level consistency. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 8866–8875.
- [39] Depu Meng, Xiaokang Chen, Zejia Fan, Gang Zeng, Houqiang Li, Yuhui Yuan, Lei Sun, and Jingdong Wang. 2021. Conditional detr for fast training convergence. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 3651–3660.
- [40] Alex Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. 2021. Learning Transferable Visual Models From Natural Language Supervision.
- [41] Ilija Radosavovic, Tete Xiao, Stephen James, Pieter Abbeel, Jitendra Malik, and Trevor Darrell. 2022. Real-World Robot Learning with Masked Visual Pre-training. *CoRL* (2022).
- [42] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. 2016. You only look once: Unified, real-time object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 779–788.
- [43] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. 2015. Faster r-cnn: Towards real-time object detection with region proposal networks. *Advances in neural information processing systems* 28 (2015).
- [44] Farzaneh Rezaeianaran, Rakshit Shetty, Rahaf Aljundi, Daniel Olmeda Reino, Shanshan Zhang, and Bernt Schiele. 2021. Seeking Similarities over Differences: Similarity-based Domain Alignment for Adaptive Object Detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 9204–9213.
- [45] Kuniaki Saito, Yoshitaka Ushiku, Tatsuya Harada, and Kate Saenko. 2019. Strong-weak distribution alignment for adaptive object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 6956–6965.
- [46] Christos Sakaridis, Dengxin Dai, and Luc Van Gool. 2018. Semantic foggy scene understanding with synthetic data. *International Journal of Computer Vision* 126, 9 (2018), 973–992.
- [47] Christos Sakaridis, Dengxin Dai, and Luc Van Gool. 2021. ACDC: The adverse conditions dataset with correspondences for semantic driving scene understanding. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*.

- 10765–10775.
- [48] Mark Sandler, Andrey Zhmoginov, Max Vladymyrov, and Andrew Jackson. 2022. Fine-tuning Image Transformers using Learnable Memory.
- [49] Max Schwarz, Anton Milan, Arul Selvam Periyasamy, and Sven Behnke. 2018. RGB-D object detection and semantic segmentation for autonomous manipulation in clutter. *The International Journal of Robotics Research* 37, 4-5 (2018), 437–451.
- [50] Zhi Tian, Chunhua Shen, Hao Chen, and Tong He. 2019. Fcos: Fully convolutional one-stage object detection. In *Proceedings of the IEEE/CVF international conference on computer vision*. 9627–9636.
- [51] Laurens Van der Maaten and Geoffrey Hinton. 2008. Visualizing data using t-SNE. *Journal of machine learning research* 9, 11 (2008).
- [52] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems* 30 (2017).
- [53] Wen Wang, Yang Cao, Jing Zhang, Fengxiang He, Zheng-Jun Zha, Yonggang Wen, and Dacheng Tao. 2021. Exploring Sequence Feature Alignment for Domain Adaptive Detection Transformers. In *Proceedings of the 29th ACM International Conference on Multimedia*. 1730–1738.
- [54] Zifeng Wang, Zizhao Zhang, Sayna Ebrahimi, Ruoxi Sun, Han Zhang, Chen-Yu Lee, Xiaoqi Ren, Guolong Su, Vincent Perot, Jennifer Dy, et al. 2022. Dualprompt: Complementary prompting for rehearsal-free continual learning. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXVI*. Springer, 631–648.
- [55] Zifeng Wang, Zizhao Zhang, Chen-Yu Lee, Han Zhang, Ruoxi Sun, Xiaoqi Ren, Guolong Su, Vincent Perot, Jennifer Dy, and Tomas Pfister. 2022. Learning to Prompt for Continual Learning. In *CVPR*.
- [56] Chang-Dong Xu, Xing-Ran Zhao, Xin Jin, and Xiu-Shen Wei. 2020. Exploring categorical regularization for domain adaptive object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 11724–11733.
- [57] Minghao Xu, Hang Wang, Bingbing Ni, Qi Tian, and Wenjun Zhang. 2020. Cross-domain detection via graph-induced prototype alignment. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 12355–12364.
- [58] Senqiao Yang, Jiarui Wu, Jiaming Liu, Xiaoqi Li, Qizhe Zhang, Mingjie Pan, and Shanghang Zhang. 2023. Exploring sparse visual prompt for cross-domain semantic segmentation. *arXiv preprint arXiv:2303.09792* (2023).
- [59] Ze Yang, Shaohui Liu, Han Hu, Liwei Wang, and Stephen Lin. 2019. Reppoints: Point set representation for object detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 9657–9666.
- [60] Zhengyuan Yang, Songyang Zhang, Liwei Wang, and Jiebo Luo. 2021. Sat: 2d semantics assisted training for 3d visual grounding. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 1856–1866.
- [61] Yuan Yao, Ao Zhang, Zhengyan Zhang, Zhiyuan Liu, Tat-Seng Chua, and Maosong Sun. 2021. CPT: Colorful Prompt Tuning for Pre-trained Vision-Language Models.
- [62] Fisher Yu, Wenqi Xian, Yingying Chen, Fangchen Liu, Mike Liao, Vashisht Madhavan, and Trevor Darrell. 2018. Bdd100k: A diverse driving video database with scalable annotation tooling. *arXiv preprint arXiv:1805.04687* 2, 5 (2018), 6.
- [63] Jinze Yu, Jiaming Liu, Xiaobao Wei, Haoyi Zhou, Yohei Nakata, Denis Gudovskiy, Tomoyuki Okuno, Jianxin Li, Kurt Keutzer, and Shanghang Zhang. 2022. MT-Trans: Cross-domain Object Detection with Mean Teacher Transformer. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part IX*. Springer, 629–645.
- [64] Wenbo Zhang, Ge-Peng Ji, Zhuo Wang, Keren Fu, and Qijun Zhao. 2021. Depth quality-inspired feature manipulation for efficient RGB-D salient object detection. In *Proceedings of the 29th ACM international conference on multimedia*. 731–740.
- [65] Liang Zheng, Liyue Shen, Lu Tian, Shengjin Wang, Jingdong Wang, and Qi Tian. 2015. Scalable person re-identification: A benchmark. In *Proceedings of the IEEE international conference on computer vision*. 1116–1124.
- [66] Yangtao Zheng, Di Huang, Songtao Liu, and Yunhong Wang. 2020. Cross-domain object detection through coarse-to-fine feature adaptation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 13766–13775.
- [67] Chaoliang Zhong, Jie Wang, Cheng Feng, Ying Zhang, Jun Sun, and Yasuto Yokota. 2022. Pica: point-wise instance and centroid alignment based few-shot domain adaptive object detection with loose annotations. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*. 2329–2338.
- [68] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. 2022. Conditional Prompt Learning for Vision-Language Models. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [69] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. 2022. Learning to Prompt for Vision-Language Models. *International Journal of Computer Vision (IJCV)* (2022).
- [70] Xinge Zhu, Jiangmiao Pang, Ceyuan Yang, Jianping Shi, and Dahua Lin. 2019. Adapting object detectors via selective cross-domain alignment. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 687–696.
- [71] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. 2020. Deformable detr: Deformable transformers for end-to-end object detection. *arXiv preprint arXiv:2010.04159* (2020).