



# Robust Test-Time Adaptation for Zero-Shot Prompt Tuning

**Ding-Chu Zhang\*, Zhi Zhou\*, Yu-Feng Li<sup>†</sup>**

National Key Laboratory for Novel Software Technology, Nanjing University, China

School of Artificial Intelligence, Nanjing University, China

{zhangdc,zhouz,liyf}@lamda.nju.edu.cn

## 1、DNN

**Problem:** the remarkable performance of DNNs heavily **relies on supervised training with large annotated datasets**, which can be a significant burden in terms of labor and computational costs.

## 2、UDA

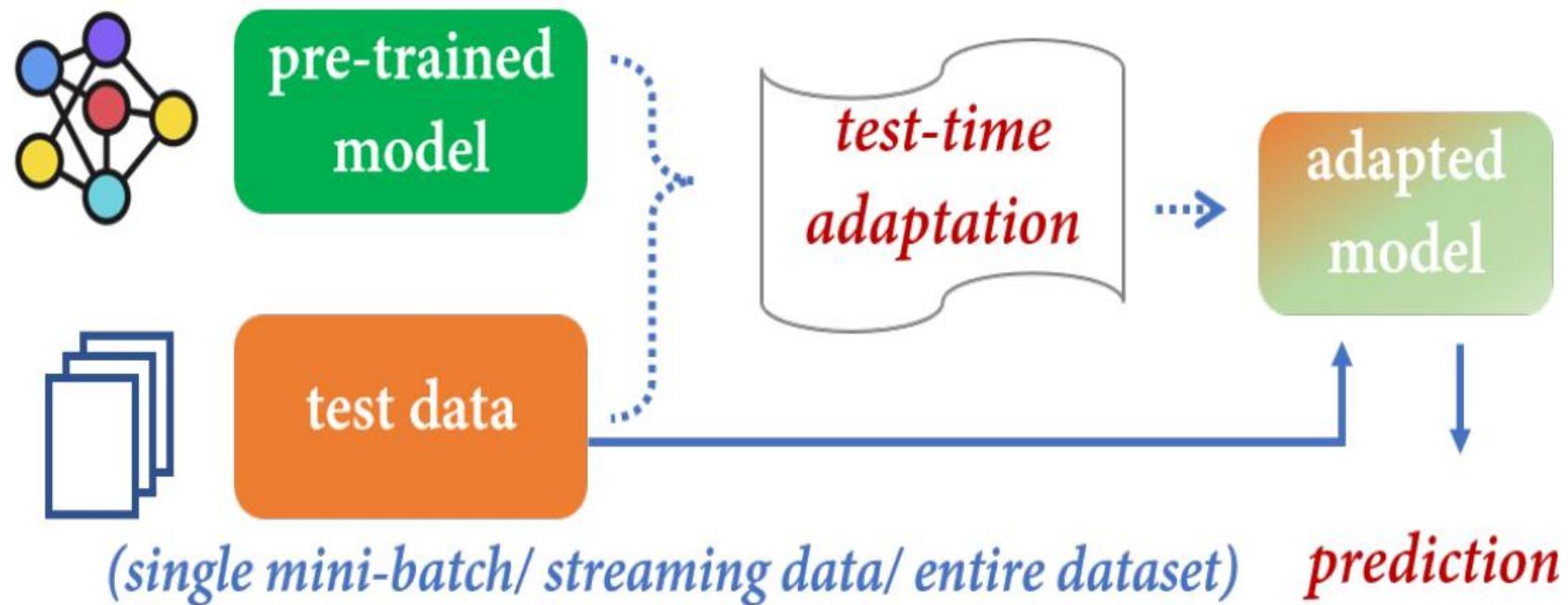
Unsupervised Domain Adaptation (UDA) is proposed to transfer knowledge from a labeled source domain to an unlabeled target domain without compromising accuracy.

**Problem:** While this approach holds promise, most UDA methods require **exposure to labeled source data** during training, which can potentially compromise **personal privacy** and **company security**.



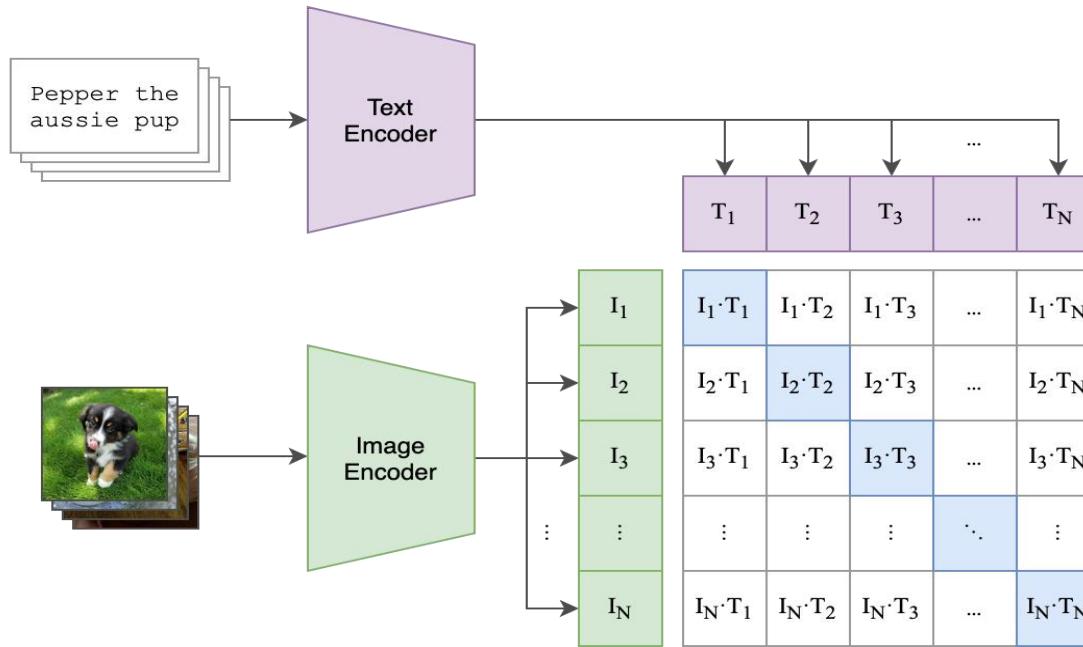
## 3、TTA

**Test-Time Domain Adaptation** (TTA) draws inspiration from test-time training and updates the model at the inference time to ensure real-time performance.

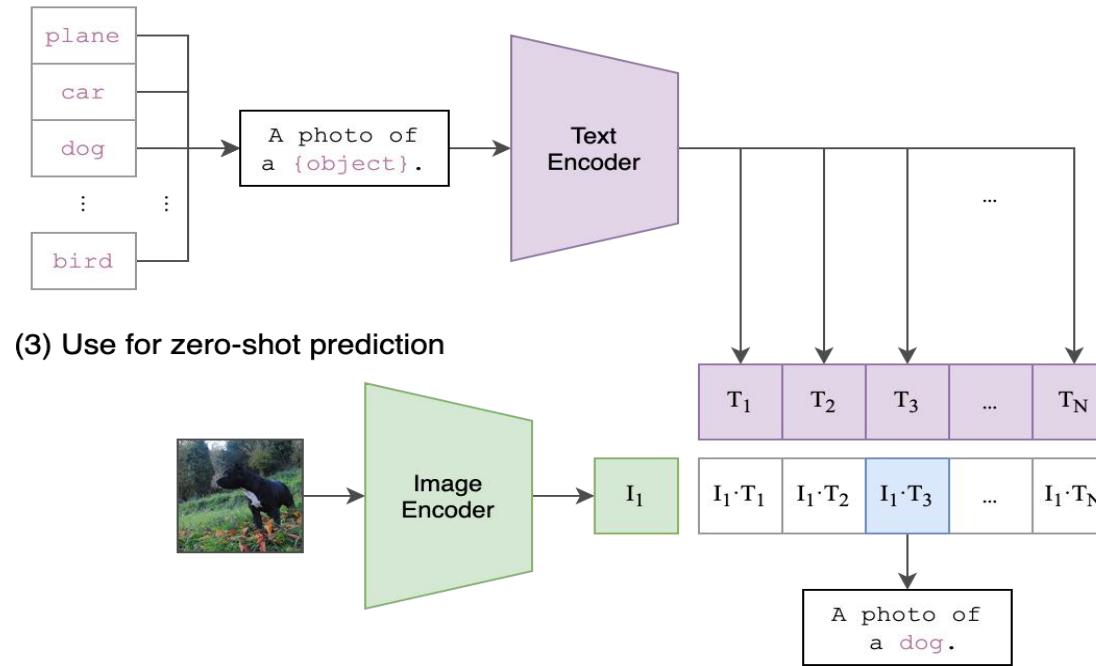


# Motivation

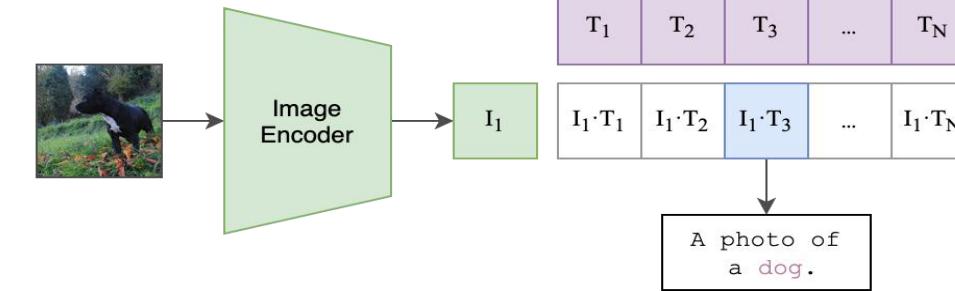
## (1) Contrastive pre-training



## (2) Create dataset classifier from label text



## (3) Use for zero-shot prediction



**Data Bias** causes a problem that **the performance of different prompts can vary across datasets**, resulting in the difficulty of selecting an optimal prompt for downstream tasks.

**Model Bias** causes prediction biases towards **specific classes, leading to error accumulation**. And the errors accumulated by Model Bias will finally result in performance degradation problem.

# Methods

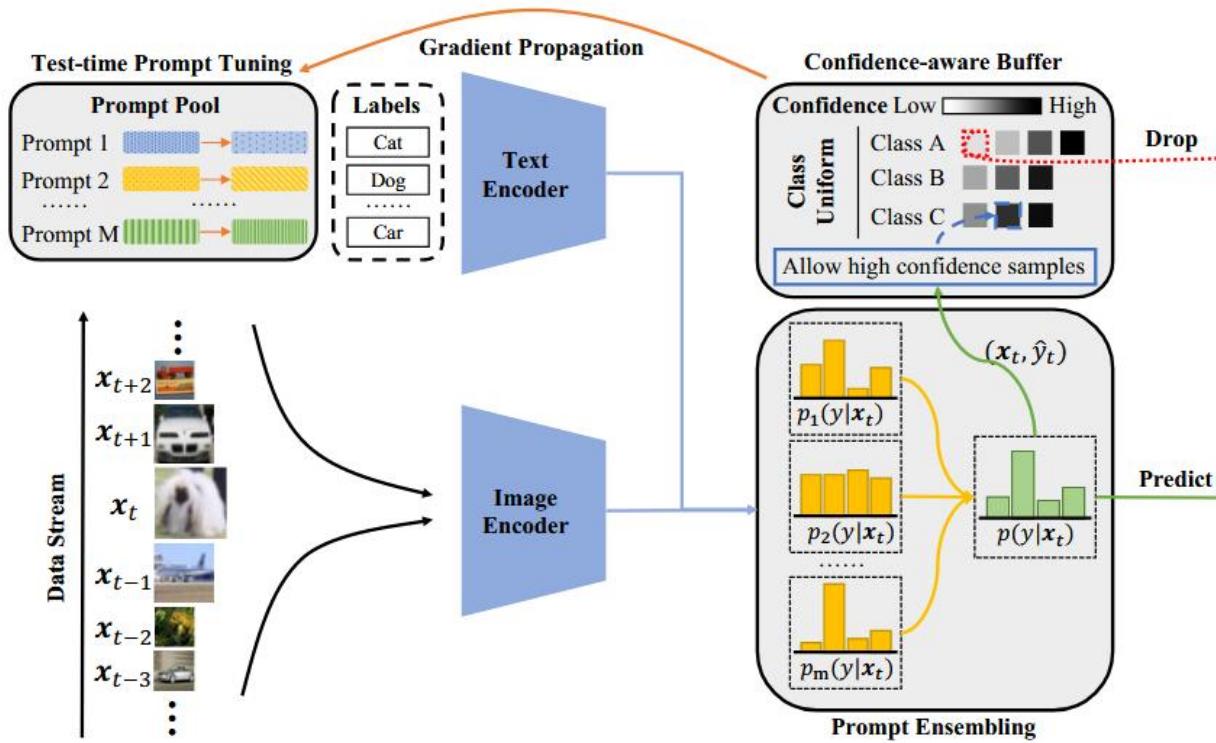


Figure 3: ADAPROMPT for image classification. We select confident samples from unlabeled test data stream by ensembling multiple prompts and push them into the confidence-aware buffer, which is used to store confident and balanced samples. Then, ADAPROMPT extracts samples from the buffer and adaptively updates all the prompts, which adapts prompts to current data.

## Data Bias

to alleviate the negative effects of *Data Bias* and avoid the worst-case results. Let  $M$  denote the number of prompts we use. We obtain the ensembled probability of different prompts in the following formulation:

$$\hat{f}(y|x_t; \mathbf{p}) = \frac{1}{M} \sum_{i=1}^M f(y|x_t; \mathbf{p}^i) \quad (3)$$

Specifically, we use a common small number of hand-crafted prompts, such as "an image of a", "a colorful image of a" and "a noisy picture of a". And we obtain pseudo label and confidence for each sample in following formulation:

$$\begin{aligned} \hat{y}(\mathbf{x}_t) &= \operatorname{argmax}_k \hat{f}(y_k|\mathbf{x}_t; \mathbf{p}) \\ c(\mathbf{x}_t) &= \max_k \hat{f}(y_k|\mathbf{x}_t; \mathbf{p}) \end{aligned} \quad (4)$$

$$L(\mathbf{x}_t) = - \sum_{k=1}^K \hat{y}_k(\mathbf{x}_t) \log \hat{f}(y_k|\mathbf{x}_t; \mathbf{p}) \quad (5)$$

where  $K$  represents the number of classes and  $\hat{y}$  represents the pseudo label obtained by Eq. (4). The purpose of minimizing cross-entropy loss is to make the model more confident in the predicted samples, which can adapt prompts to *Data Bias* and improve the accuracy of predictions.

# Model Bias



- 当 updateCriterion 为 random 时：

- 如果要移除的实例所属类别不在存储数据量最大的类别索引列表（largest\_indices）中，那么从最大的类别里随机选择一个（通过 `random.choice`），再在选中的最大类别对应的数据列表里随机选择一个索引（通过 `random.randrange`）作为要移除的目标索引 `tgt_idx`，然后从该类别对应的各个维度数据列表（特征、类别、其他信息）中移除对应索引的数据。
- 如果要移除的实例所属类别就在最大类别索引列表中，会根据当前类别已存储的数据量 `m_c`、该类别总的实际计数 `n_c` 以及一个随机数 `u` 来判断是否移除，若  $u \leq m_c / n_c$ ，则同样随机选择一个索引 `tgt_idx` 并移除对应的数据，否则返回 `False`，表示不移除。

- 当 updateCriterion 为 FIFO 时：

- 如果要移除的实例所属类别不在最大类别索引列表中，直接选择最大类别对应数据列表中的第一个索引（即按照先进先出原则）作为 `tgt_idx`，然后移除对应的数据。
- 如果要移除的实例所属类别在最大类别索引列表中，同样选择该类别对应数据列表中的第一个索引进行移除操作。

- 当 updateCriterion 为 Threshold 时：

- 如果要移除的实例所属类别不在最大类别索引列表中，选择最大类别对应存储的置信度列表中最小值对应的索引（即选择置信度最小的实例对应的索引）作为 `tgt_idx`，然后移除对应的数据。
- 如果要移除的实例所属类别在最大类别索引列表中，先获取该类别对应置信度列表中的最小值 `minThreshold`，如果传入的实例置信度小于这个最小值，则返回 `False`，不移除；否则选择最小值对应的索引 `tgt_idx` 并移除对应的数据。

# Experiments



Dataset		CIFAR10-C(s=3)			CIFAR10-C(s=5)			CIFAR100-C(s=3)			CIFAR100-C(s=5)		
Methods		Source	TPT	Ours	Source	TPT	Ours	Source	TPT	Ours	Source	TPT	Ours
Noise	Gauss.	50.03	52.86	<b>54.50</b>	38.00	40.08	<b>42.48</b>	27.81	25.54	<b>28.61</b>	19.60	17.31	<b>21.92</b>
	Shot	61.74	63.32	<b>64.92</b>	43.14	44.74	<b>47.89</b>	33.81	32.22	<b>35.30</b>	21.36	19.04	<b>23.95</b>
	Impul.	78.59	78.87	<b>81.36</b>	56.70	59.08	<b>60.59</b>	47.30	47.63	<b>50.51</b>	25.31	25.65	<b>30.06</b>
Blur	Defoc.	85.46	85.25	<b>87.69</b>	72.88	72.10	<b>74.98</b>	60.10	<b>60.55</b>	60.54	42.52	42.73	<b>43.07</b>
	Glass	54.26	53.95	<b>59.29</b>	42.59	43.19	<b>47.51</b>	29.35	29.21	<b>30.38</b>	20.06	19.97	<b>20.91</b>
	Motion	77.15	77.06	<b>78.52</b>	70.96	70.14	<b>72.54</b>	48.69	48.86	<b>49.69</b>	<b>43.15</b>	42.63	42.46
	Zoom	81.57	81.35	<b>84.29</b>	74.66	74.89	<b>78.30</b>	56.08	55.96	<b>57.22</b>	47.89	48.12	<b>48.72</b>
Weather	Snow.	81.01	81.18	<b>84.52</b>	74.74	75.32	<b>78.26</b>	53.90	55.41	<b>56.34</b>	48.35	<b>49.19</b>	48.95
	Frost	81.13	81.02	<b>84.60</b>	78.40	78.33	<b>80.19</b>	53.12	53.89	<b>55.05</b>	49.72	50.43	<b>50.89</b>
	Fog	86.60	86.49	<b>89.10</b>	71.66	72.54	<b>73.14</b>	60.77	<b>61.64</b>	61.33	41.64	<b>42.71</b>	42.45
	Brit.	88.92	88.67	<b>91.53</b>	85.00	85.12	<b>88.06</b>	64.88	65.39	<b>66.64</b>	57.02	57.58	<b>59.07</b>
Digital	Contr.	87.11	87.70	<b>89.28</b>	63.00	<b>70.80</b>	67.95	59.77	61.18	<b>61.58</b>	34.54	<b>38.06</b>	36.84
	Elastic	80.27	80.75	<b>83.46</b>	55.40	<b>57.10</b>	<b>58.88</b>	52.53	53.43	<b>55.01</b>	29.21	30.05	<b>30.56</b>
	Pixel	75.18	75.98	<b>81.54</b>	48.09	52.24	<b>57.21</b>	51.09	51.94	<b>53.29</b>	23.94	25.15	<b>27.50</b>
	JPEG	69.51	69.82	<b>72.67</b>	60.30	61.55	<b>63.83</b>	39.68	40.17	<b>42.40</b>	32.46	32.43	<b>34.29</b>
Avg.		75.90	76.29	<b>79.15</b>	62.37	63.81	<b>66.12</b>	49.26	49.54	<b>50.93</b>	35.78	36.07	<b>37.44</b>

Table 1: Comparison with state-of-the-art test-time prompt tuning methods on CIFAR10-C and CIFAR100-C benchmarks with corruption level 3 and 5. We conduct separate tests on 15 different domains for each benchmark. We omit std in this table due to space issues. The best results are indicated in bold. Our method outperforms comparison methods in almost all cases. The best performance is in bold.

# Experiments

Methods		Source	TPT	TPT-C	Ours
Noise	Gauss.	15.72	16.29	0.52	<b>17.52</b>
	Shot	23.44	23.86	0.52	<b>26.47</b>
	Impul.	17.47	17.58	0.52	<b>20.76</b>
Blur	Defoc.	32.43	32.65	0.58	<b>34.39</b>
	Glass	11.88	12.51	0.52	<b>14.45</b>
	Motion	31.97	32.31	0.54	<b>33.98</b>
	Zoom	30.99	31.57	0.54	<b>33.32</b>
Weather	Snow.	29.69	30.90	0.55	<b>32.82</b>
	Frost	32.98	33.25	0.58	<b>36.30</b>
	Fog	35.81	36.36	0.58	<b>37.97</b>
	Brit.	43.95	43.62	0.60	<b>46.80</b>
Digital	Contr.	22.56	23.00	0.52	<b>25.52</b>
	Elastic	38.14	38.74	0.58	<b>40.78</b>
	Pixel	26.38	27.72	0.55	<b>29.42</b>
	JPEG	37.54	37.56	0.64	<b>40.72</b>
Avg.		28.73	29.20	0.55	<b>31.42</b>

Table 3: Comparison with SOTA test-time prompt tuning methods on TinyImageNet-C with corruption level 3. ADA\_PROMPT outperforms them in all domains.

# Experiments

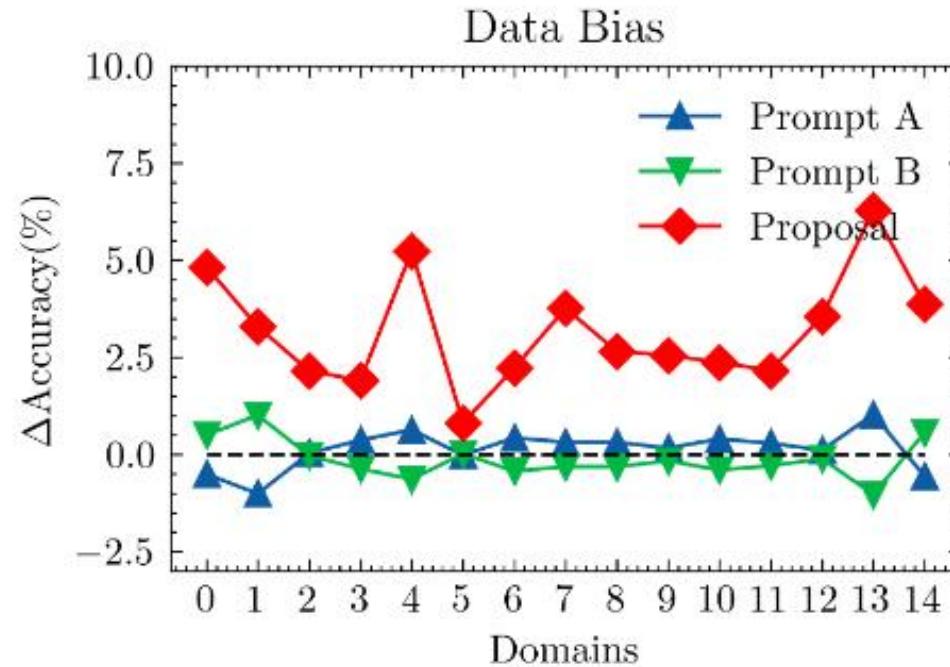


Figure 1: Relative performance compared to the average performance of prompts evaluated on CIFAR10-C in 15 domains with corruption level 3 using different initial prompts.

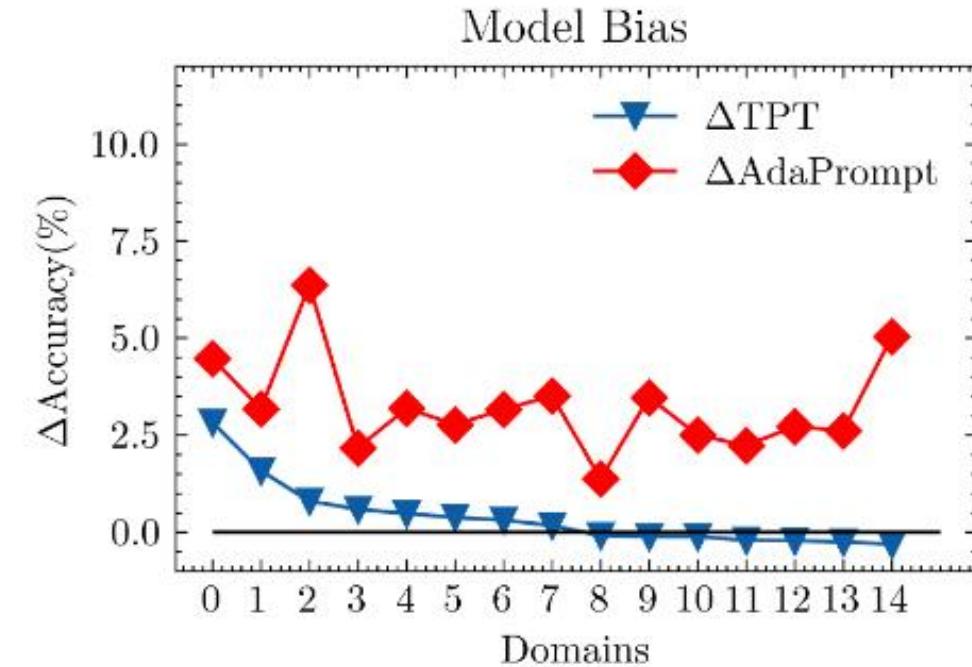


Figure 2: Relative performance compared to baseline evaluated on CIFAR10-C in 15 different domains with corruption level 3. The black line represents the baseline.

# Ablation Experiments

Method	CIFAR10-C(s=3)	CIFAR10-C(s=5)
$P_A$	$75.91 \pm 0.00$	$62.37 \pm 0.00$
$P_B$	$76.21 \pm 0.00$	$62.77 \pm 0.00$
$P_C$	$72.98 \pm 0.00$	$59.25 \pm 0.00$
$P_{best} + UP.$	$77.72 \pm 0.24$	$65.32 \pm 0.18$
$P_e$	$75.38 \pm 0.00$	$61.75 \pm 0.00$
$P_e + UP.$	<b><math>79.15 \pm 0.23</math></b>	<b><math>66.12 \pm 0.43</math></b>

Table 2: Evaluation of each module on CIFAR10-C with corruption level 3 and 5. The average accuracy of different modules on 15 different domains is shown.

Component	$M_e$	$M_u$	CIFAR10-C(s=3)	CIFAR10-C(s=5)
		✓	$76.21 \pm 0.00$	$62.37 \pm 0.00$
		✓	$75.38 \pm 0.00$	$61.75 \pm 0.00$
	✓	✓	$77.72 \pm 0.24$	$65.32 \pm 0.18$
	✓	✓	<b><math>79.15 \pm 0.23</math></b>	<b><math>66.12 \pm 0.43</math></b>

Table 4: Ablation study of ADAPROMPT on CIFAR10-C dataset with corruption level 3 and 5. The average accuracy on 15 different domains is reported.

Acc(%)	Source	TPT	Ours
RN50	$47.70 \pm 0.00$	$51.44 \pm 0.02$	<b><math>55.44 \pm 0.30</math></b>
ViT-B/32	$71.30 \pm 0.00$	$73.77 \pm 0.03$	<b><math>75.81 \pm 0.33</math></b>

Table 5: Average accuracy of CIFAR10-C in different 15 domains with corruption level 3 on different backbones.

# Ablation Experiments

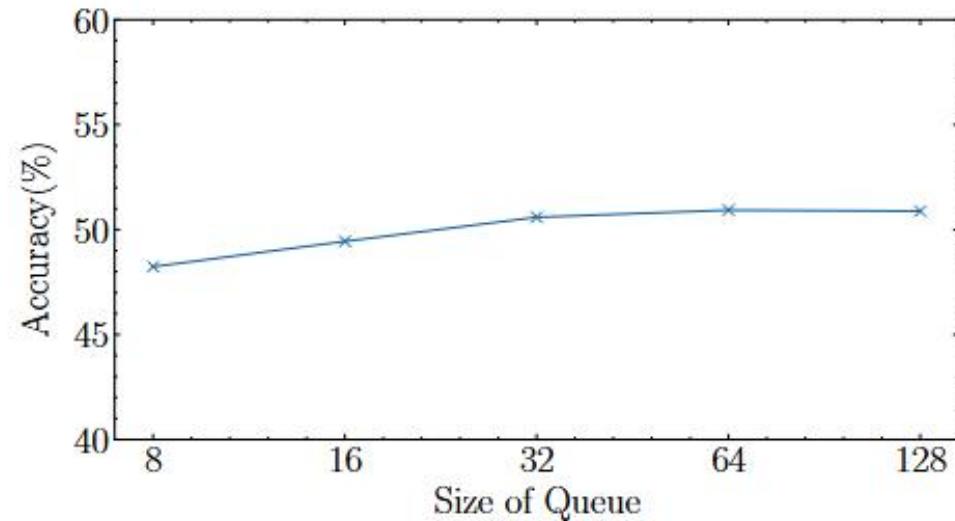


Figure 5: Average accuracy of ADA\_PROMPT with different buffer size on CIFAR100-C with corruption level 3.

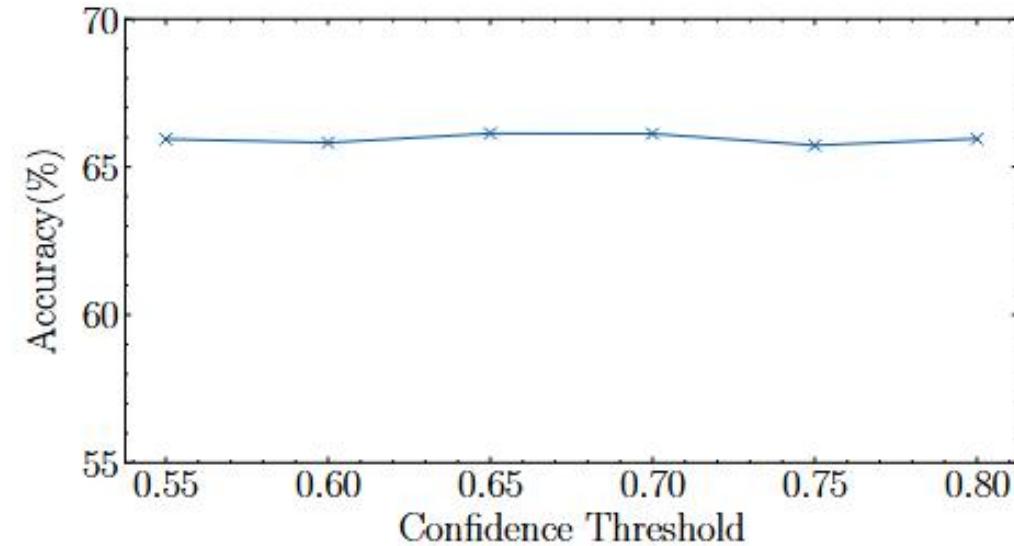


Figure 6: Performance of ADA\_PROMPT with different confidence threshold on CIFAR100-C with corruption level 3.

Thanks