

# Efficient Test-Time Adaptation

Related works and Project Status

Juan Camacho, Andrea De Carlo, Samuele Bolotta

# Motivation

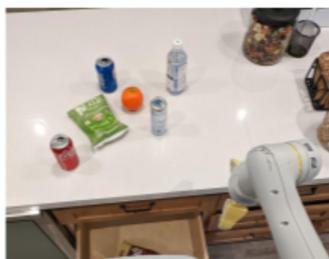
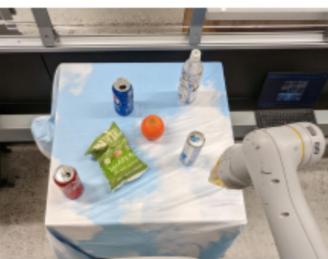


RT-1: Robotics Transformer for Real-World Control at Scale, Google Research, 2022

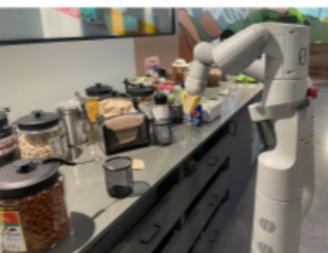
# Motivation



Robustness to the number of distractor objects



robustness to different backgrounds and environments such as new, previously unseen kitchens



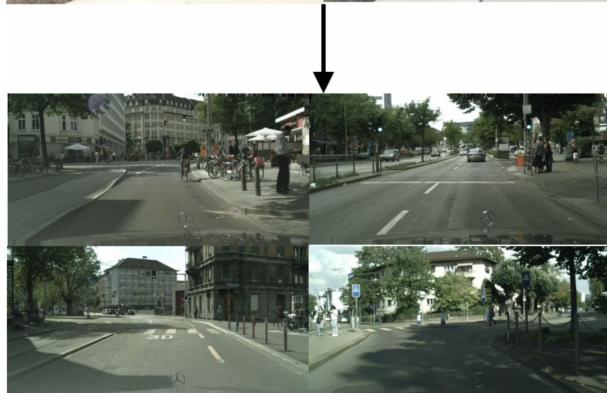
realistic scenarios that combine all these elements

# Motivation

Sim to real



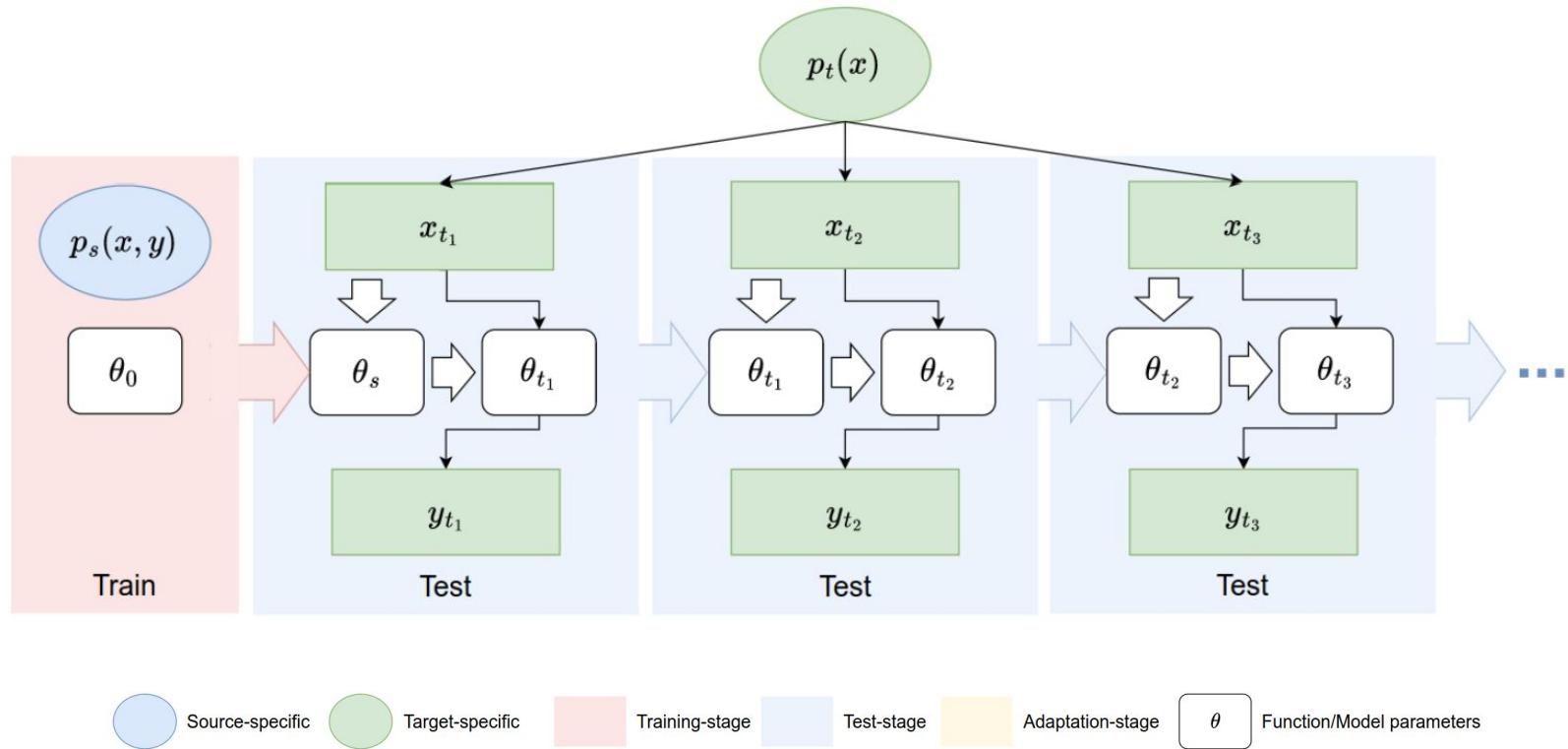
Different conditions



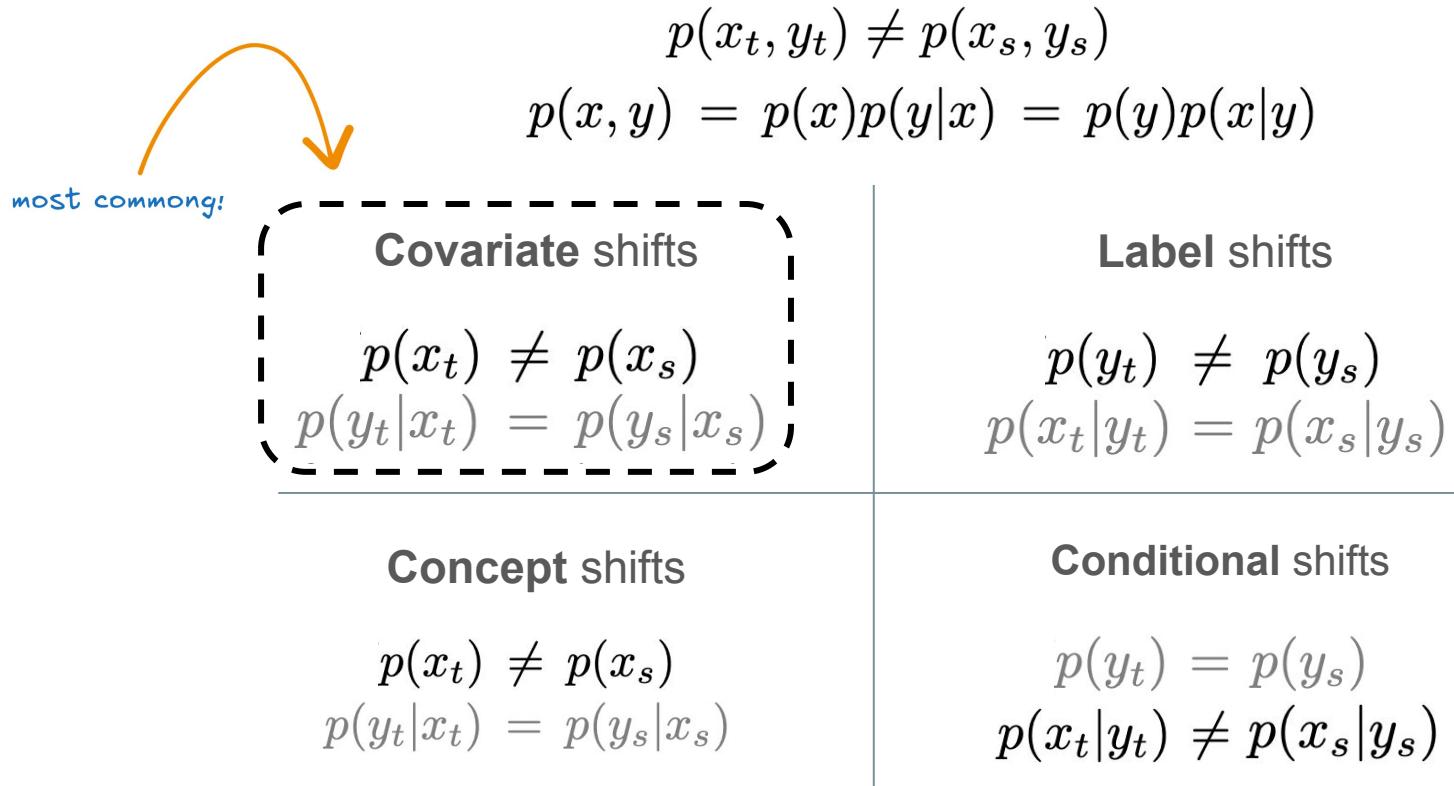
# Motivation



# But what is Test-Time Adaptation (TTA)?

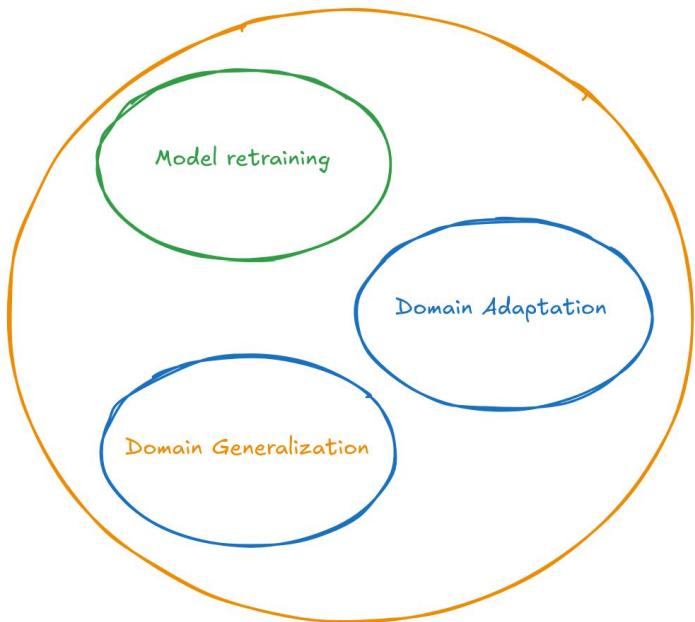


# TTA addresses distribution shifts



# Why Test-Time Adaptation?

## Traditional solutions

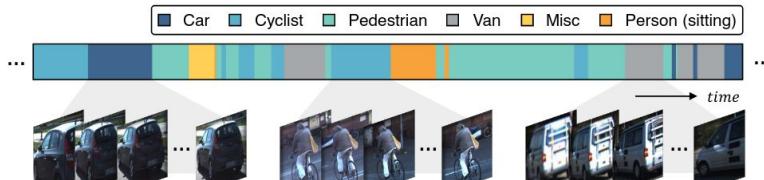


## TTA advantages

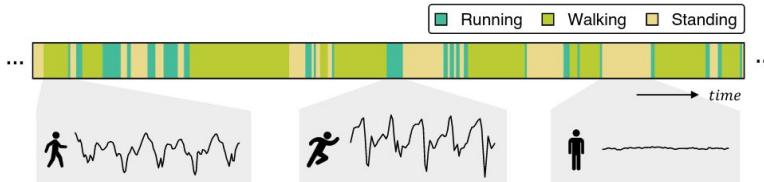
- Online Adaptation
- Resource efficiency
- Practical benefits

# Understanding data streams

## Non-I.I.D. Data Stream

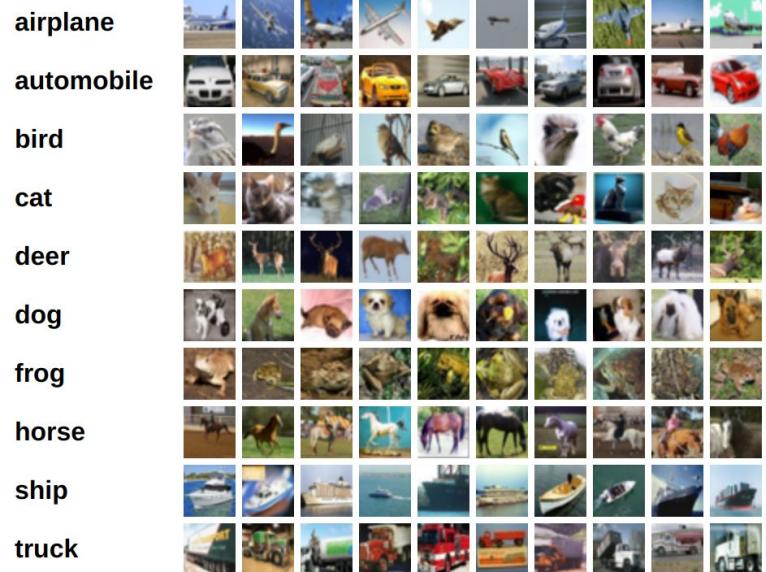


(a) KITTI dataset.



(b) HARTH dataset.

## I.I.D. Data Stream



CIFAR-10

NOTE: Robust Continual Test-time Adaptation Against Temporal Correlation

# Understanding data streams

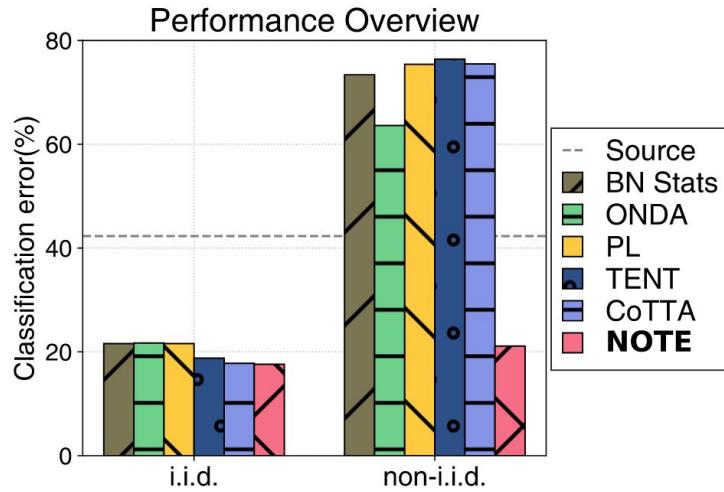


Figure 2: Average classification error (%) of existing TTA methods and our method (NOTE) on CIFAR10-C [13]. The error rates significantly increase under the non-i.i.d. setting compared with the i.i.d. setting. Lower is better.

# Impact of Data Streams on TTA

## I.I.D. Data Streams

- Simpler adaptation methods
- Efficiency
- Fast performance

## Non-I.I.D. Data Streams

- Complex adaptation strategies
- Temporal aware
- Resource heavy

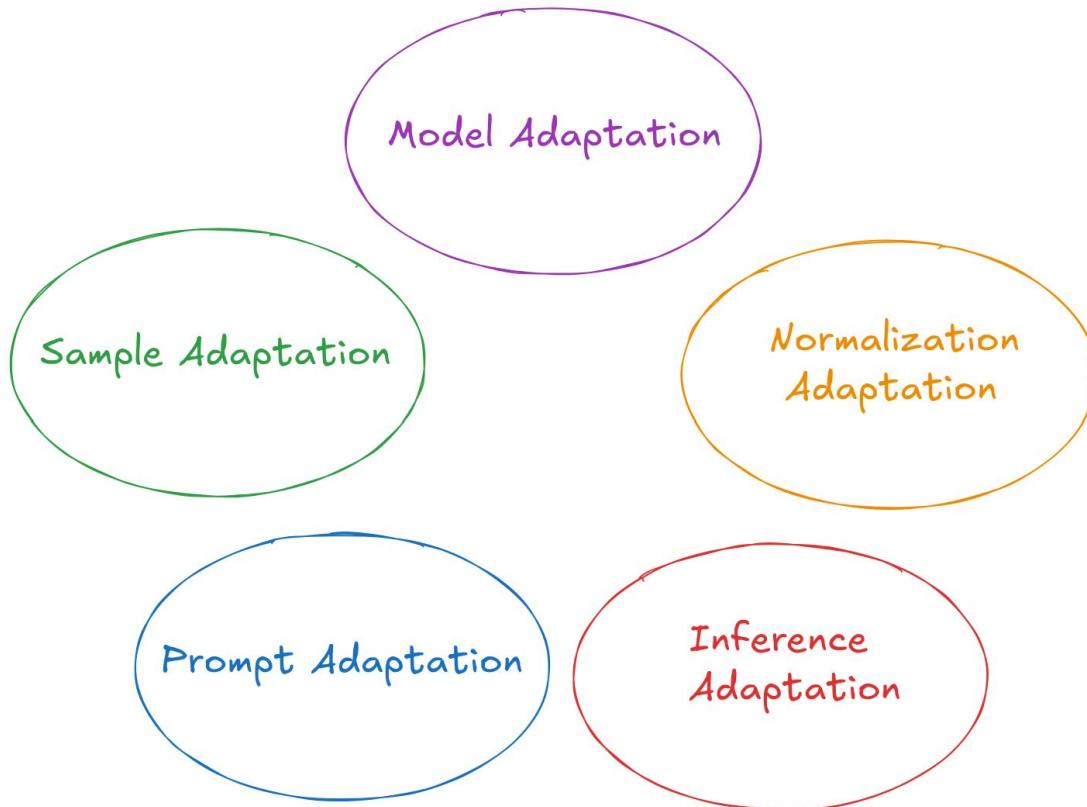
## Key Trade-offs

 Accuracy vs Speed

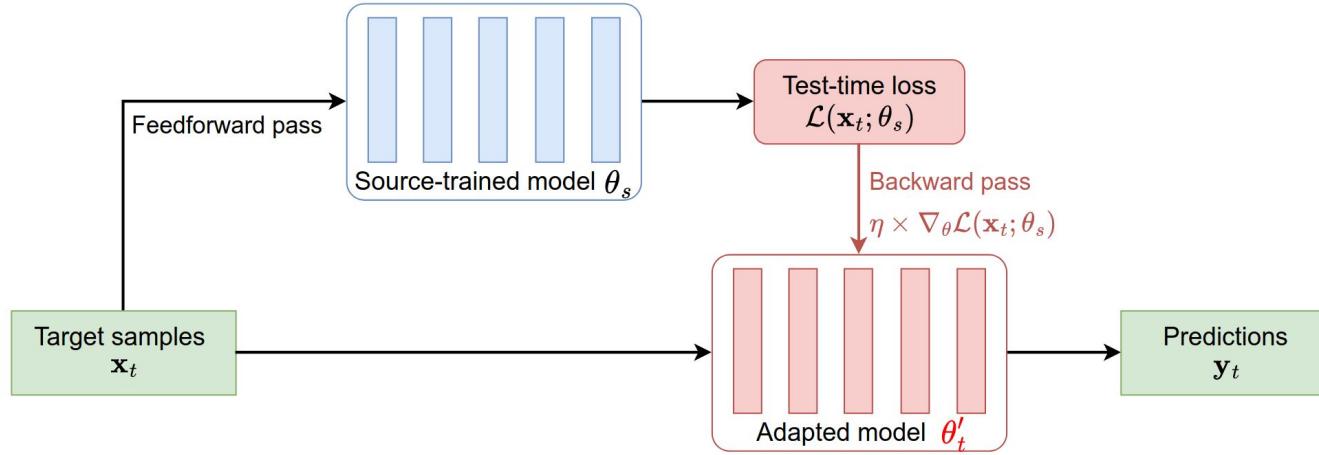
 Stability vs Plasticity

 Complexity vs Efficiency

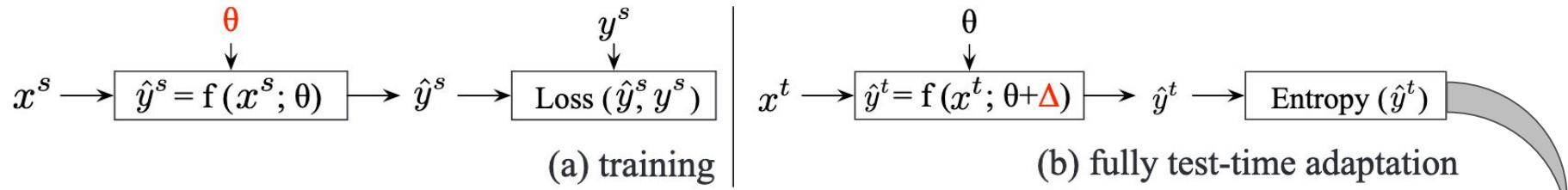
# Test-Time Adaptation Methods



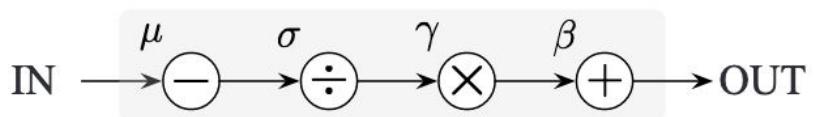
# Model Adaptation



# Model Adaptation

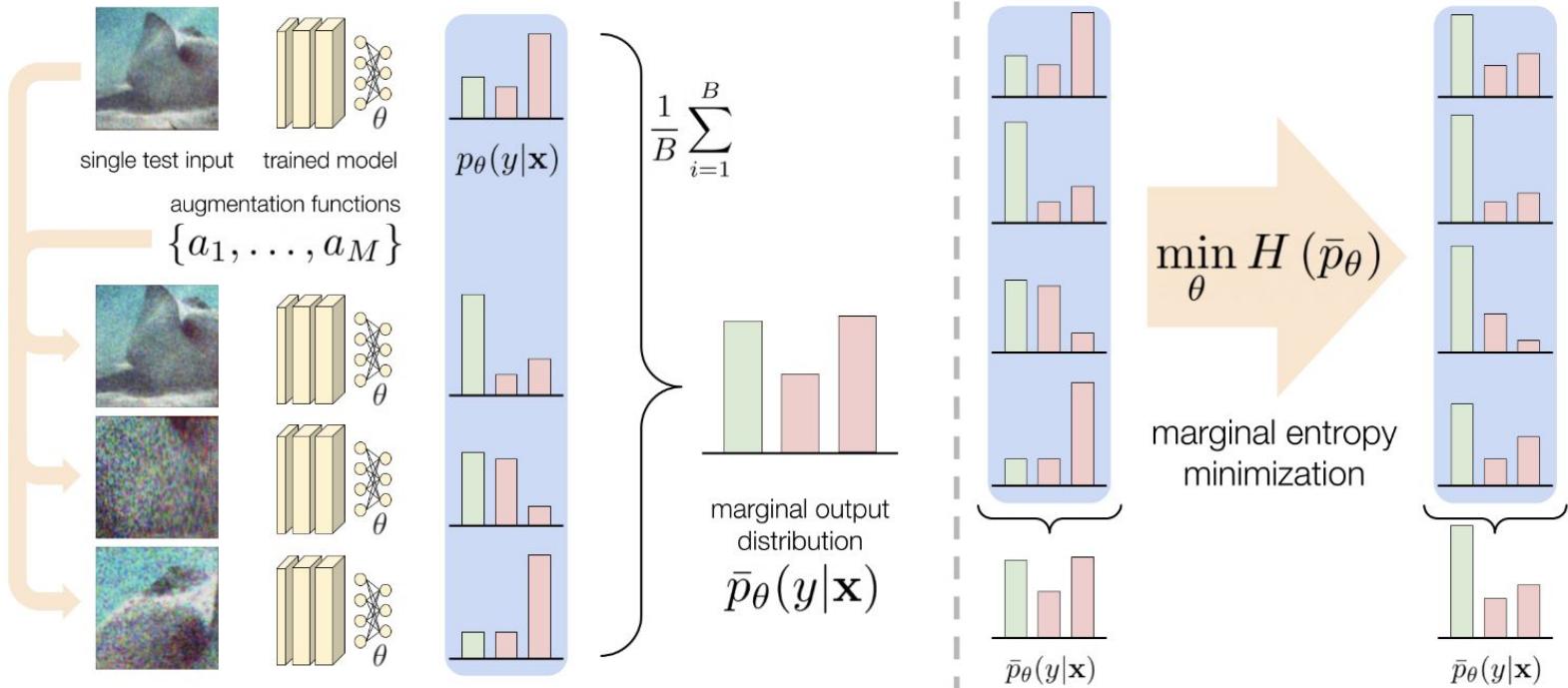


$$H(\hat{y}) = - \sum_c p(\hat{y}_c) \log p(\hat{y}_c)$$



normalization  $\mu \leftarrow \mathbb{E}[x_t], \sigma^2 \leftarrow \mathbb{E}[(\mu - x_t)^2]$   
transformation  $\gamma \leftarrow \gamma + \partial H / \partial \gamma, \beta \leftarrow \beta + \partial H / \partial \beta$

# Model Adaptation



# Model Adaptation

Test-time adaptation with “zero” temperature

Marginal probability distribution:

$$\bar{p}(\cdot|\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N p(\cdot|\mathbf{x}_i)$$

Marginal entropy loss:

$$\mathcal{L}_{ent} = H(\bar{p}(\cdot|\mathbf{x})) = - \sum_{c=1}^C \bar{p}(y = c|\mathbf{x}) \log(\bar{p}(y = c|\mathbf{x}))$$

# Model Adaptation

Test-time adaptation with “zero” temperature

Marginal entropy loss:

$$\mathcal{L}_{ent} = H(\bar{p}(\cdot|\mathbf{x})) = - \sum_{c=1}^C \bar{p}(y = c|\mathbf{x}) \log(\bar{p}(y = c|\mathbf{x}))$$

Re-written loss in TPT:

$$\mathcal{L}_{ent} = H(\bar{p}(\cdot|\mathbf{x}, \mathbf{t}_{ctx})) = - \sum_c^C \bar{p}(y = c|\mathbf{x}, \mathbf{t}_{ctx}, \tau) \log(\bar{p}(y = c|\mathbf{x}, \mathbf{t}_{ctx}, \tau))$$

$$\text{where } \bar{p}(y = c|\mathbf{x}, \mathbf{t}_{ctx}, \tau) = \frac{1}{N} \sum_i^N \frac{\exp\left(\mathbf{z}_i^{\text{img}} \cdot \mathbf{z}_c^{\text{txt}}(\mathbf{t}_{ctx})/\tau\right)}{\sum_k^C \exp\left(\mathbf{z}_i^{\text{img}} \cdot \mathbf{z}_k^{\text{txt}}(\mathbf{t}_{ctx})/\tau\right)}.$$

# Model Adaptation

Test-time adaptation with zero temperature

```
def zero(image, z_txt, N, gamma, temp):
    """
    :param z_txt: pre-computed text embeddings (C,hdim)
    :param temp: model's original temperature
    :param augment: takes (C,H,W) and returns (N,C,H,W)
    :param gamma: filtering percentile (e.g., 0.3)
    """
    views = augment(image, num_views=N) # generate augmented views
    l = model.image_encoder(views) @ z_txt.t() # predict (unscaled logits)
    l_filt = confidence_filter(l, temp, top=gamma) # retain most confident preds
    zero_temp = torch.finfo(l_filt.dtype).eps # zero temperature
    p_bar = (l_filt / zero_temp).softmax(dim=1).sum(dim=0) # marginalize
    return p_bar.argmax()
```

# Model Adaptation

## Test-time adaptation with zero temperature

Table 2: Fine-grained classification. TTA methods are grouped according to the reference baseline, top-1 accuracy is reported and **bold text** indicates the best performer of each group.

Method	FLWR	DTD	PETS	CARS	UCF	CAL	FOOD	SUN	AIR	ESAT	Mean	Median
CLIP-ViT-B-16												
Zero-Shot	67.44	44.27	88.25	65.48	65.13	93.35	83.65	62.59	23.67	42.01	63.58	65.31
Ensemble	67.07	45.09	<b>88.28</b>	66.16	67.51	93.91	84.04	66.26	23.22	<b>50.42</b>	65.20	66.66
TPT	<b>68.75</b>	<b>47.04</b>	87.23	66.68	68.16	93.93	84.67	65.39	23.13	42.86	64.78	67.42
ZERO	67.68	46.12	87.75	68.04	67.77	93.66	86.53	65.03	<b>25.21</b>	34.33	64.21	67.72
ZERO+Ensemble	67.17	45.86	87.83	<b>68.97</b>	<b>69.18</b>	<b>94.41</b>	<b>86.77</b>	<b>67.63</b>	<b>25.21</b>	42.17	<b>65.52</b>	<b>68.30</b>
MaPLe												
Zero-Shot	72.23	46.49	90.49	65.57	68.69	93.53	86.20	67.01	24.74	<b>48.06</b>	66.30	67.85
TPT	72.37	45.87	90.72	66.50	69.19	93.59	86.64	67.54	24.70	47.80	66.49	68.36
PromptAlign	<b>72.39</b>	47.24	<b>90.76</b>	68.50	69.47	94.01	86.65	67.54	24.80	47.86	<b>66.92</b>	68.99
ZERO	71.62	<b>47.89</b>	90.60	<b>68.58</b>	<b>69.87</b>	<b>94.48</b>	<b>87.20</b>	<b>68.20</b>	<b>26.25</b>	39.47	66.42	<b>69.23</b>
CLIP-ViT-B-16 + CLIP-ViT-L-14												
Zero-Shot	75.76	51.83	92.86	76.16	73.70	94.04	88.03	66.96	30.54	<b>54.38</b>	70.43	74.73
RLCF $t_{ctx}^{t=1}$	71.58	50.34	89.01	69.76	69.84	94.09	85.90	67.33	23.71	46.87	66.84	69.80
RLCF $t_{ctx}^{t=3}$	72.49	51.93	89.55	72.91	72.31	95.00	86.84	69.04	25.40	45.96	68.14	72.40
RLCF $\Theta_v^{t=1}$	72.56	52.21	89.51	63.12	71.49	94.65	86.90	68.50	24.06	47.74	67.07	70.00
RLCF $\Theta_v^{t=3}$	71.74	53.27	91.15	70.93	73.24	94.73	87.28	69.38	28.54	47.41	68.77	71.34
ZERO	<b>76.41</b>	<b>53.63</b>	<b>94.08</b>	<b>78.39</b>	<b>74.68</b>	<b>95.21</b>	<b>90.66</b>	<b>69.61</b>	<b>33.62</b>	44.21	<b>71.05</b>	<b>75.55</b>

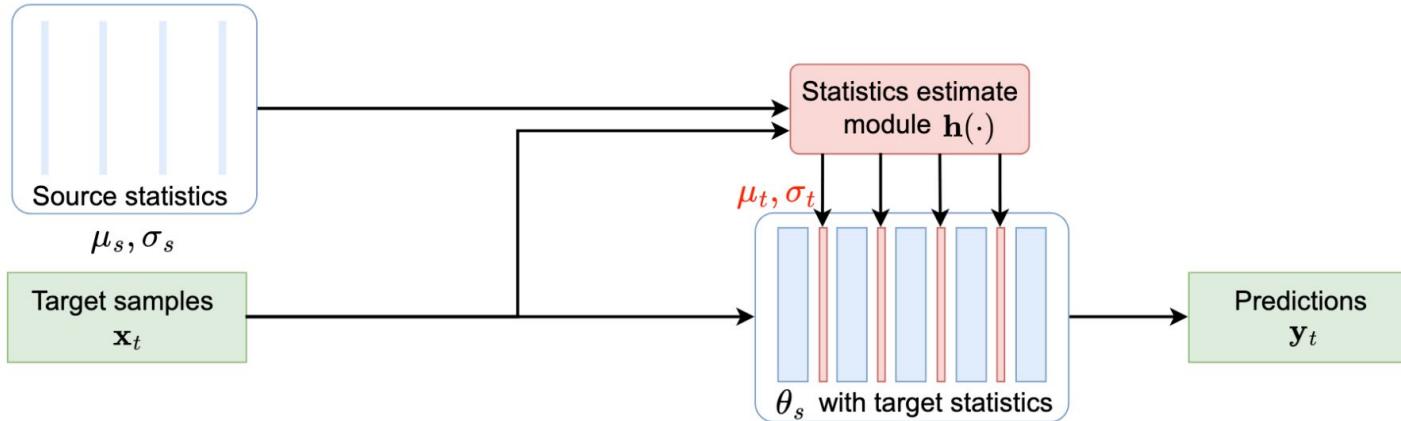
# Model Adaptation

Test-time adaptation with zero temperature

Table 3: Computational requirements of different TTA methods.

Metric	CLIP-ViT-B-16		CLIP-ViT-B-16 + CLIP-ViT-L-14		
	TPT	ZERO	RLCF $t_{\text{ctx}}=3$	RLCF $\Theta_v$ $t=3$	ZERO
Time [s]	0.57±0.01	<b>0.06±0.01</b>	1.20±0.02	0.18±0.01	<b>0.08±0.02</b>
Mem [GB]	17.66	<b>1.40</b>	18.64	9.04	<b>2.58</b>

# Normalization Adaptation



$$\hat{\mathbf{x}}_s = \frac{\mathbf{x}_s - \mu_s}{\sqrt{\sigma_s^2 + \epsilon}} \longrightarrow \mu_t, \sigma_t = \mathbf{h}(\mathbf{x}_t, \mu_s, \sigma_s), \quad \hat{\mathbf{x}}_t = \frac{\mathbf{x}_t - \mu_t}{\sqrt{\sigma_t^2 + \epsilon}}, \quad \mathbf{y}_t = f_{\theta_s}(\hat{\mathbf{x}}_t),$$

# Normalization Adaptation

$$\mu_t, \sigma_t = \mathbf{h}(\mathbf{x}_t, \mu_s, \sigma_s), \quad \hat{\mathbf{x}}_t = \frac{\mathbf{x}_t - \mu_t}{\sqrt{\sigma_t^2 + \epsilon}}, \quad \mathbf{y}_t = f_{\theta_s}(\hat{\mathbf{x}}_t),$$

## Target statistics

- Directly estimate
- Alleviate covariate shift

## Statistics combination

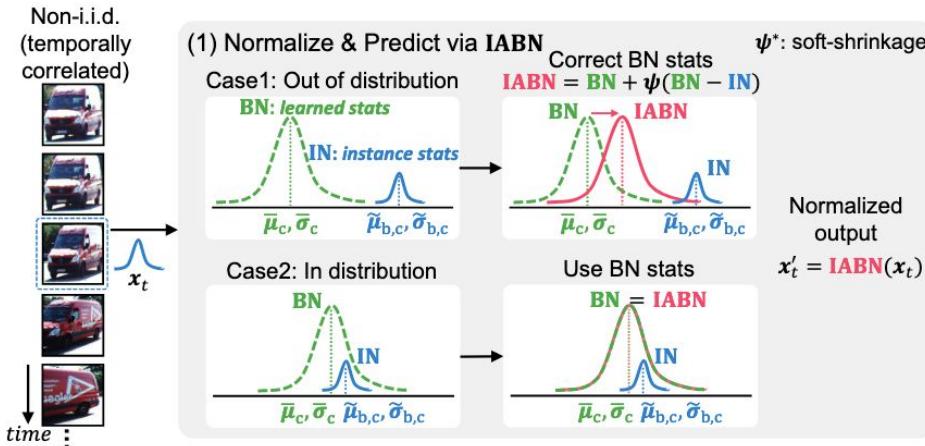
- Weighed sum
- adaptation  preservation

## Statistics inference

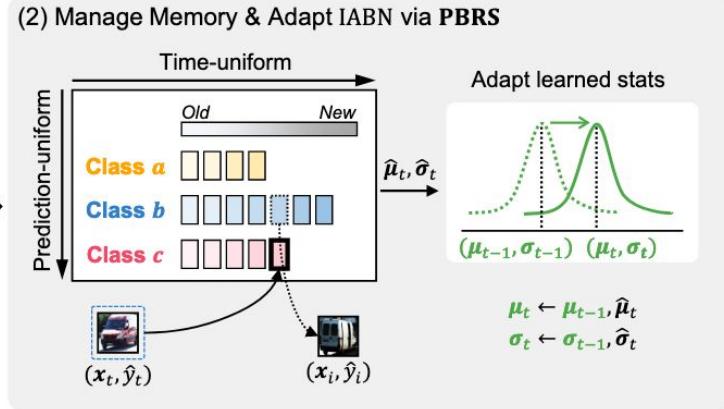
- Meta Learning per sample
- No requirements on data

# Normalization Adaptation

## Instance-Aware Batch Normalization

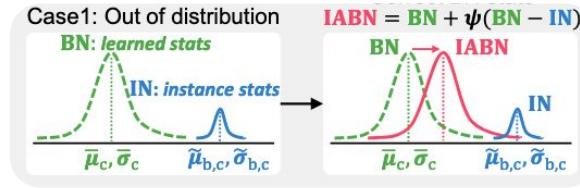
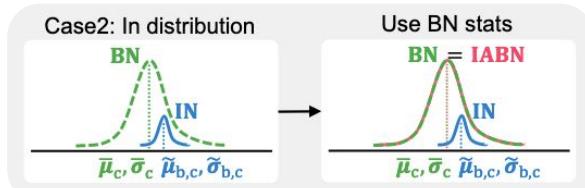


## Prediction Balanced Reservoir Sampling



# Normalization Adaptation

## Instance-Aware Batch Normalization



$$BN(\mathbf{f}_{:,c,:}; \boldsymbol{\mu}_c, \boldsymbol{\sigma}_c^2) := \gamma \cdot \frac{\mathbf{f}_{:,c,:} - \boldsymbol{\mu}_c}{\sqrt{\boldsymbol{\sigma}_c^2 + \epsilon}} + \beta$$

$$\hat{\boldsymbol{\mu}}_c := \frac{1}{BL} \sum_{b,l} \mathbf{f}_{b,c,l}$$

$$\hat{\boldsymbol{\sigma}}_c^2 := \frac{1}{BL} \sum_{b,l} (\mathbf{f}_{b,c,l} - \hat{\boldsymbol{\mu}}_c)^2$$

$$IABN(\mathbf{f}_{b,c,:}; \bar{\boldsymbol{\mu}}_c, \bar{\boldsymbol{\sigma}}_c^2; \tilde{\boldsymbol{\mu}}_{b,c}, \tilde{\boldsymbol{\sigma}}_{b,c}^2) := \gamma \cdot \frac{\mathbf{f}_{b,c,:} - \boldsymbol{\mu}_{b,c}^{IABN}}{\sqrt{(\boldsymbol{\sigma}_{b,c}^{IABN})^2 + \epsilon}} + \beta$$

$$\boldsymbol{\mu}_{b,c}^{IABN} := \bar{\boldsymbol{\mu}}_c + \psi(\tilde{\boldsymbol{\mu}}_{b,c} - \bar{\boldsymbol{\mu}}_c; \alpha s_{\tilde{\boldsymbol{\mu}},c})$$

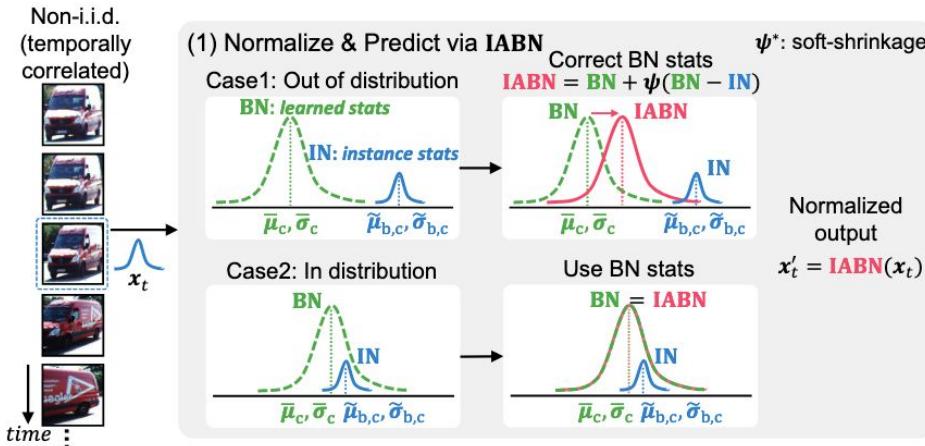
$$s_{\boldsymbol{\mu},c}^2 := \frac{\bar{\boldsymbol{\sigma}}_c^2}{L}$$

$$(\boldsymbol{\sigma}_{b,c}^{IABN})^2 := \bar{\boldsymbol{\sigma}}_c^2 + \psi(\tilde{\boldsymbol{\sigma}}_{b,c}^2 - \bar{\boldsymbol{\sigma}}_c^2; \alpha s_{\tilde{\boldsymbol{\sigma}}^2,c})$$

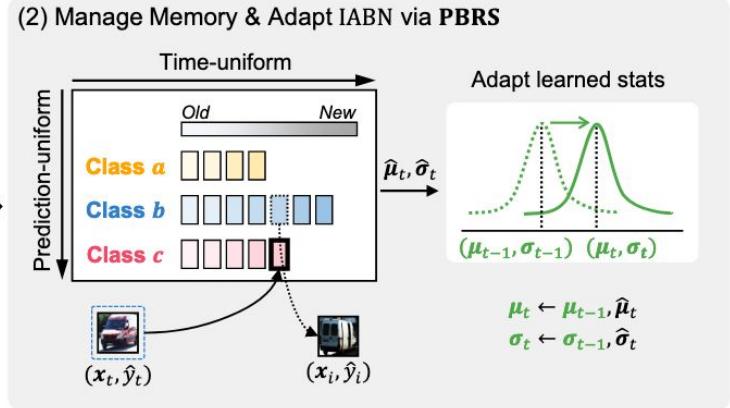
$$s_{\tilde{\boldsymbol{\sigma}}^2,c}^2 := \frac{2\bar{\boldsymbol{\sigma}}_c^4}{L-1}$$

# Normalization Adaptation

## Instance-Aware Batch Normalization



## Prediction Balanced Reservoir Sampling



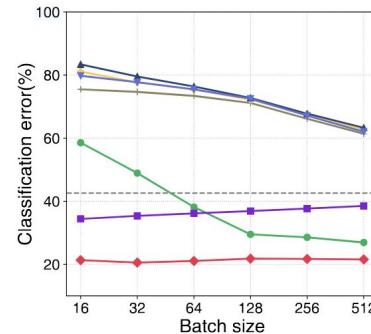
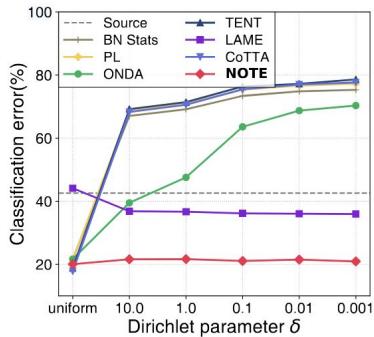
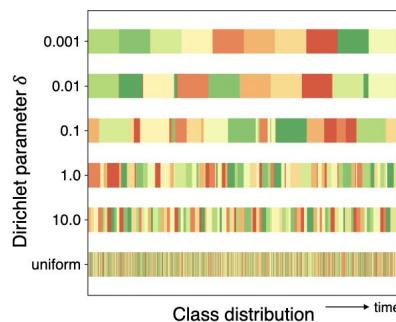
$$IABN(f_{b,c,:}; \bar{\mu}_c, \bar{\sigma}_c^2; \tilde{\mu}_{b,c}, \tilde{\sigma}_{b,c}^2) := \gamma \cdot \frac{f_{b,c,:} - \mu_{b,c}^{IABN}}{\sqrt{(\sigma_{b,c}^{IABN})^2 + \epsilon}} + \beta$$

$$\mu_t = (1 - m)\mu_{t-1} + m \frac{N}{N-1} \hat{\mu}_t$$

$$\sigma_t^2 = (1 - m)\sigma_{t-1}^2 + m \frac{N}{N-1} \hat{\sigma}_t^2$$

# Normalization Adaptation

Method	Temporally correlated test stream				Uniformly distributed test stream			
	CIFAR10-C	CIFAR100-C	ImageNet-C	Avg	CIFAR10-C	CIFAR100-C	ImageNet-C	Avg
Source	42.3 $\pm$ 1.1	66.6 $\pm$ 0.1	86.1 $\pm$ 0.0	65.0	42.3 $\pm$ 1.1	66.6 $\pm$ 0.1	86.1 $\pm$ 0.0	65.0
BN Stats [29]	73.4 $\pm$ 1.3	65.0 $\pm$ 0.3	96.9 $\pm$ 0.0	78.5	21.6 $\pm$ 0.4	46.6 $\pm$ 0.2	76.0 $\pm$ 0.0	48.1
ONDA [27]	63.6 $\pm$ 1.0	49.6 $\pm$ 0.3	89.0 $\pm$ 0.0	67.4	21.7 $\pm$ 0.4	46.5 $\pm$ 0.1	75.9 $\pm$ 0.0	48.0
PL [22]	75.4 $\pm$ 1.8	66.4 $\pm$ 0.4	98.9 $\pm$ 0.0	80.2	21.6 $\pm$ 0.2	43.1 $\pm$ 0.3	74.4 $\pm$ 0.2	46.4
TENT [41]	76.4 $\pm$ 2.7	66.9 $\pm$ 0.6	96.9 $\pm$ 0.0	80.1	18.8 $\pm$ 0.2	40.3 $\pm$ 0.2	76.0 $\pm$ 0.0	45.0
LAME [4]	36.2 $\pm$ 1.3	63.3 $\pm$ 0.3	82.7 $\pm$ 0.0	60.7	44.1 $\pm$ 0.5	68.8 $\pm$ 0.1	86.3 $\pm$ 0.0	66.4
CoTTA [44]	75.5 $\pm$ 0.7	64.2 $\pm$ 0.2	97.0 $\pm$ 0.0	78.9	17.8 $\pm$ 0.3	44.3 $\pm$ 0.2	71.5 $\pm$ 0.0	44.6
NOTE	<b>21.1 <math>\pm</math> 0.6</b>	<b>47.0 <math>\pm</math> 0.1</b>	<b>80.6 <math>\pm</math> 0.1</b>	<b>49.6</b>	20.1 $\pm$ 0.5 (17.6 $\pm$ 0.3)	46.4 $\pm$ 0.0 (41.0 $\pm$ 0.2)	<b>70.3 <math>\pm</math> 0.0</b> (71.7 $\pm$ 0.0)	45.6 (43.4)



# Normalization Adaptation

$$\mu_t, \sigma_t = \mathbf{h}(\mathbf{x}_t, \mu_s, \sigma_s), \quad \hat{\mathbf{x}}_t = \frac{\mathbf{x}_t - \mu_t}{\sqrt{\sigma_t^2 + \epsilon}}, \quad \mathbf{y}_t = f_{\theta_s}(\hat{\mathbf{x}}_t),$$

## Target statistics

- Directly estimate
- Alleviate covariate shift

Large distribution shifts  
→ mismatch

## Statistics combination

- Weighed sum
- adaptation  preservation

Hyperparameters

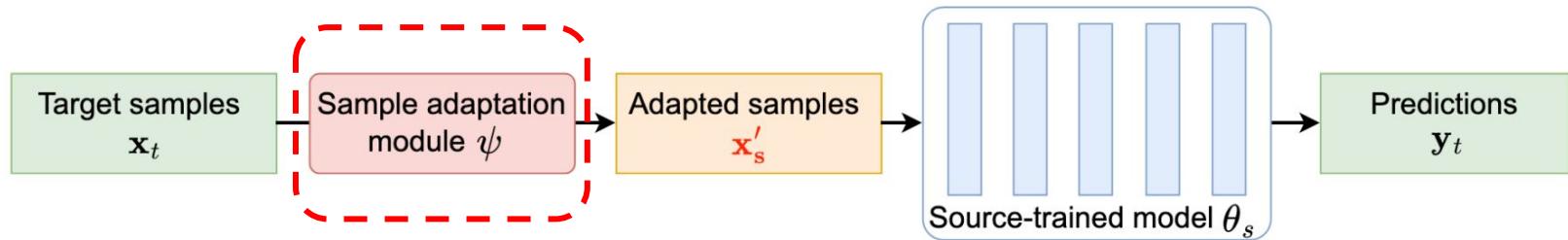
## Statistics inference

- Meta Learning per sample
- No requirements on data

Difficult to learn

NOT model-agnostic

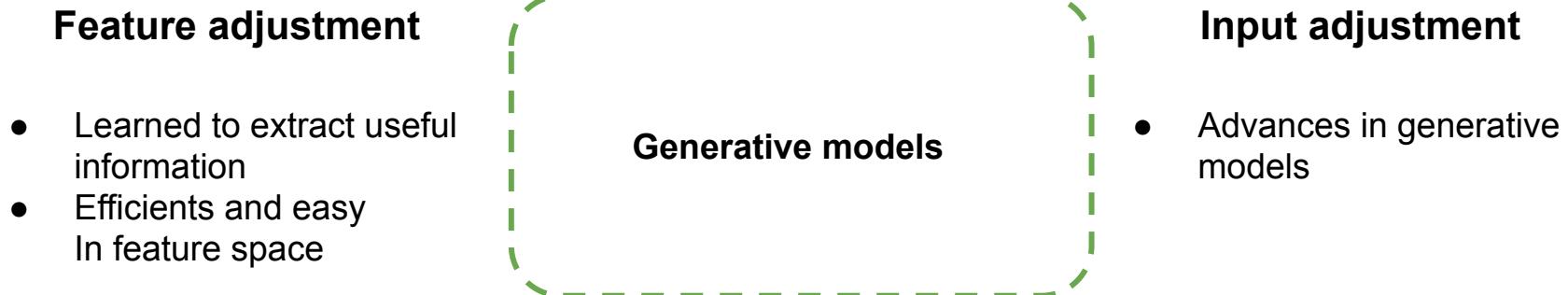
# Sample Adaptation



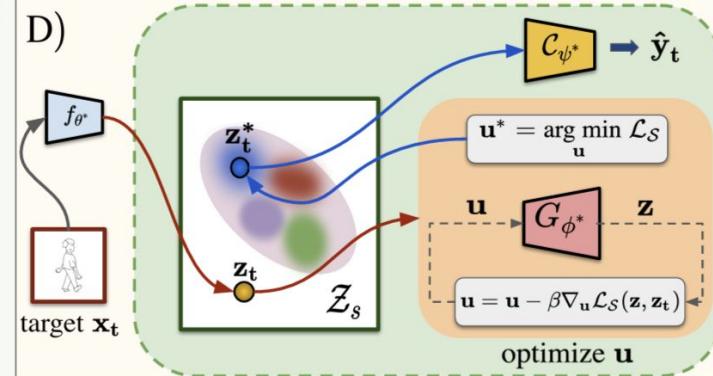
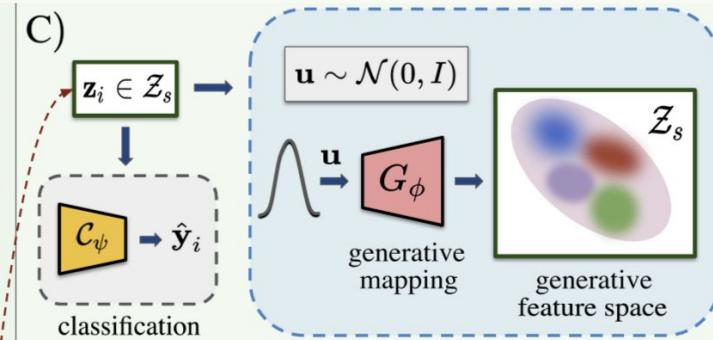
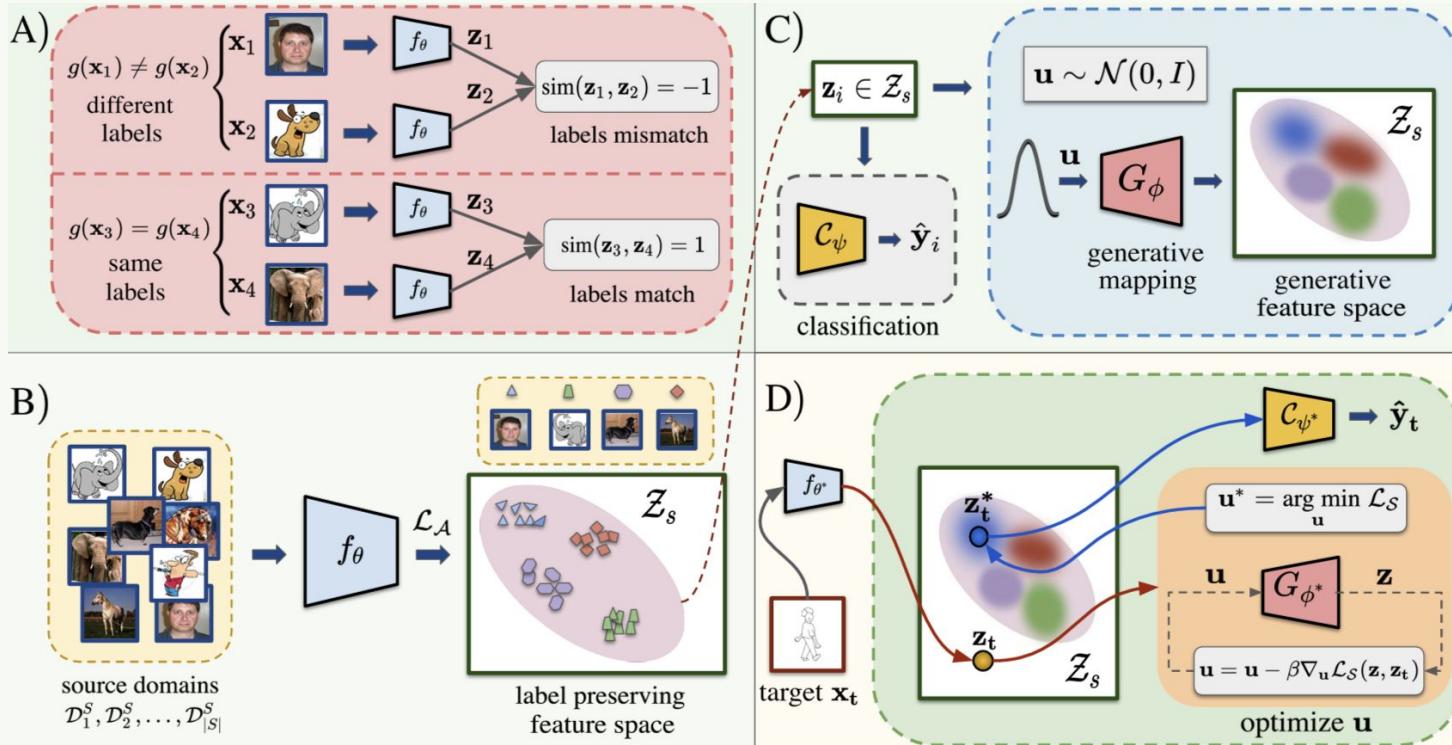
$$\mathbf{x}'_s = \psi(\mathbf{x}_t), \quad \mathbf{y}_t = f_{\theta_s}(\mathbf{x}'_s).$$

# Sample Adaptation

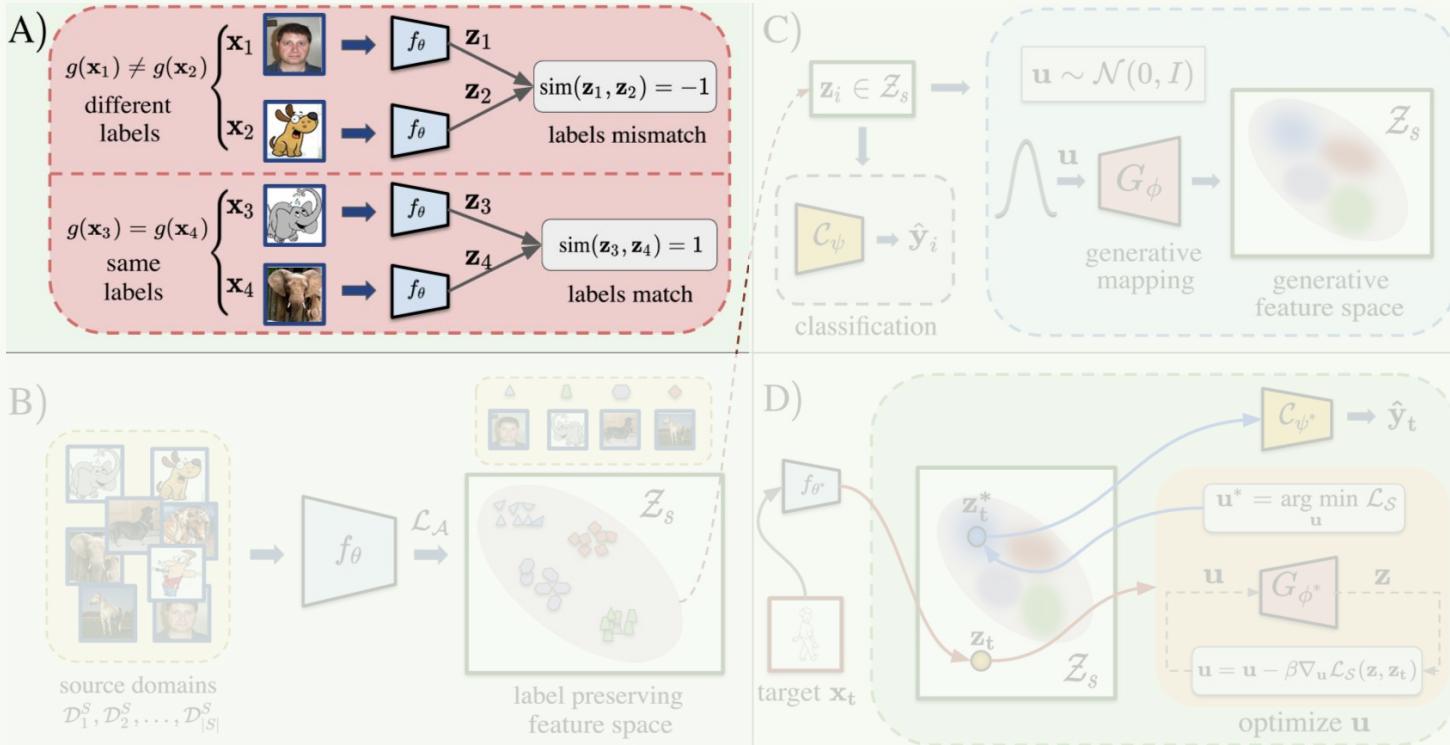
$$\mathbf{x}'_s = \psi(\mathbf{x}_t), \quad \mathbf{y}_t = f_{\theta_s}(\mathbf{x}'_s).$$



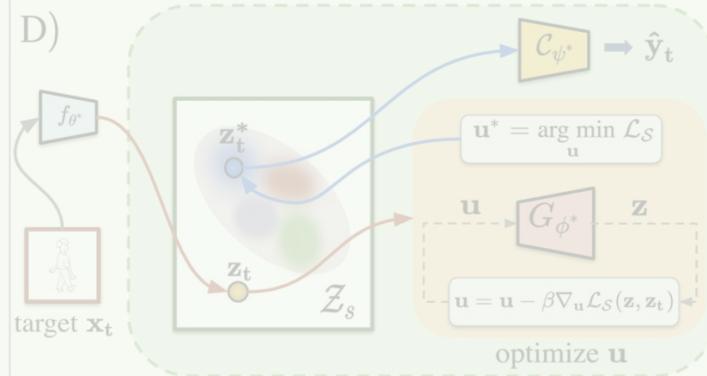
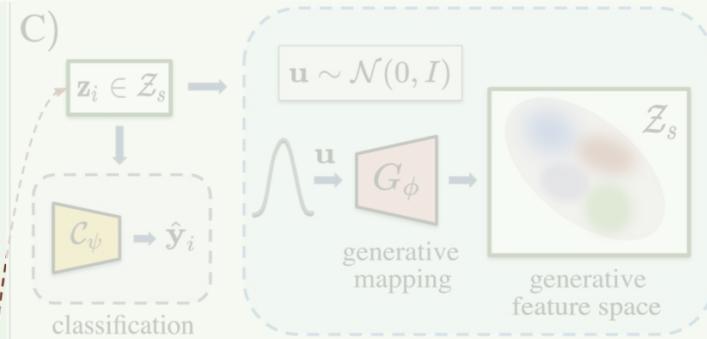
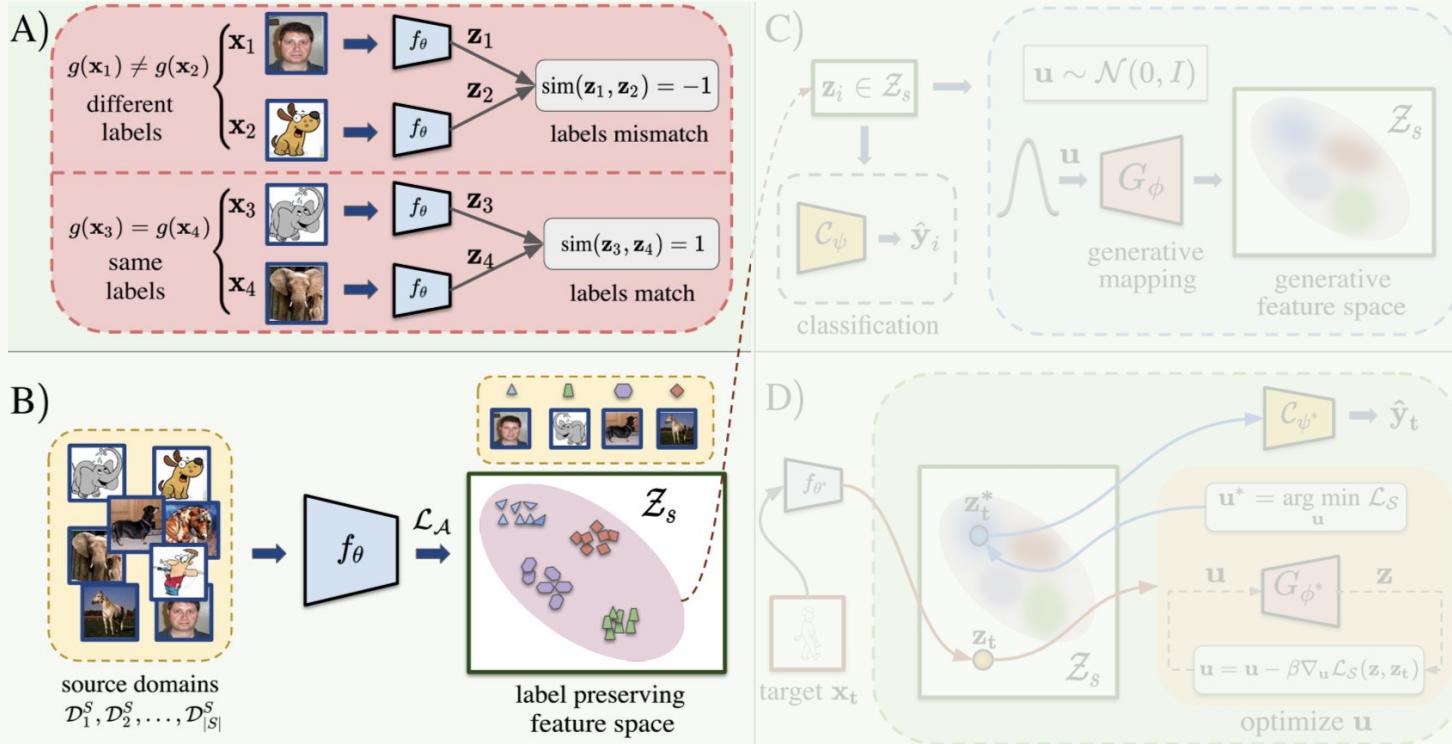
# Sample Adaptation



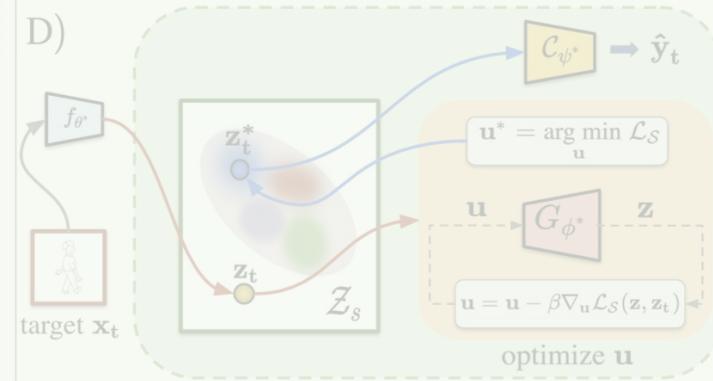
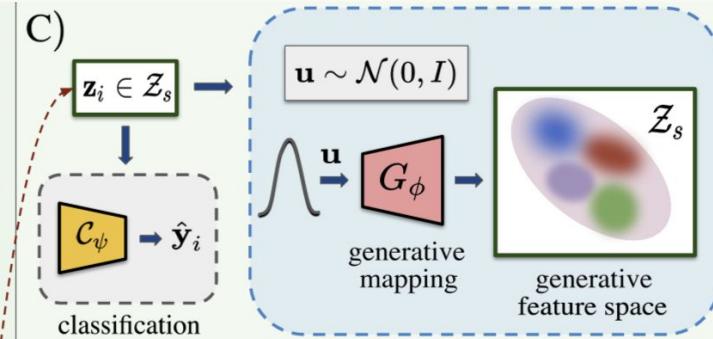
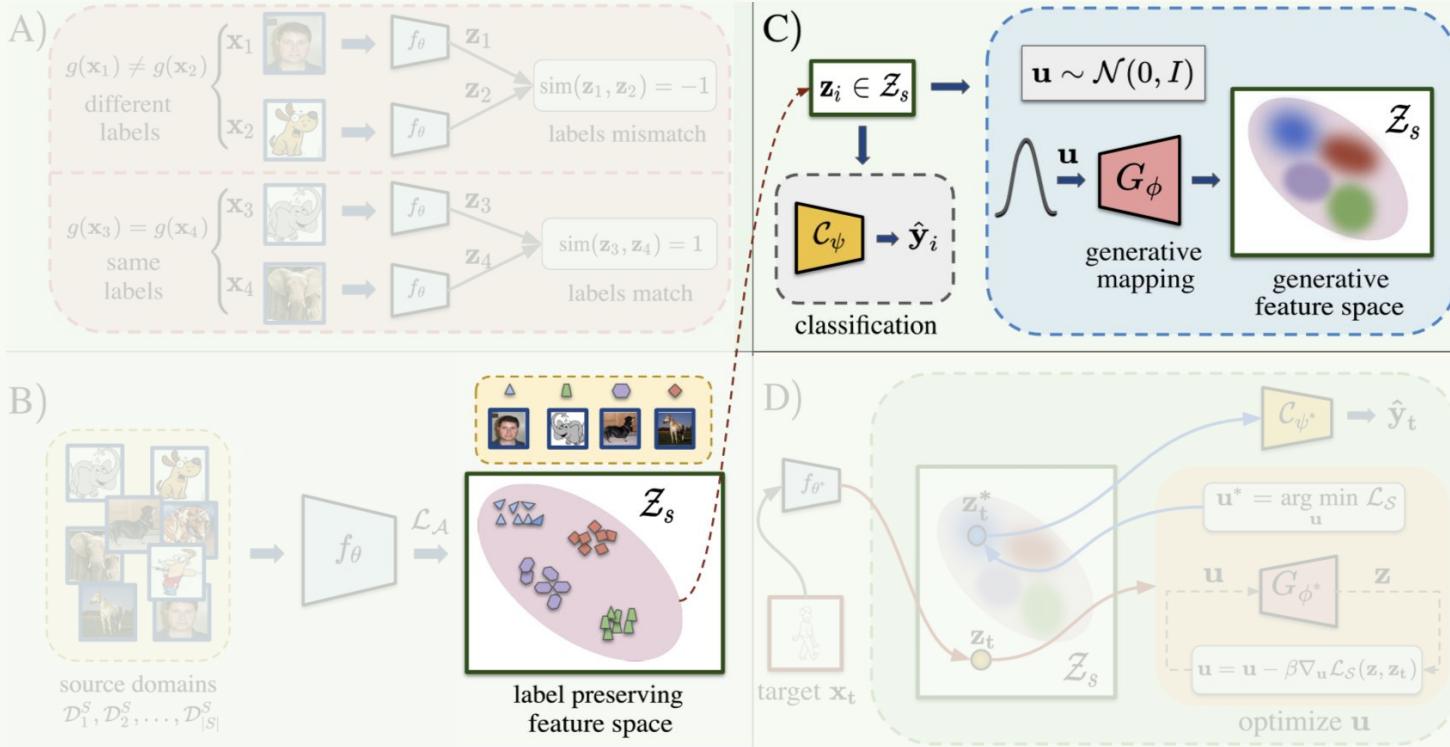
# Sample Adaptation



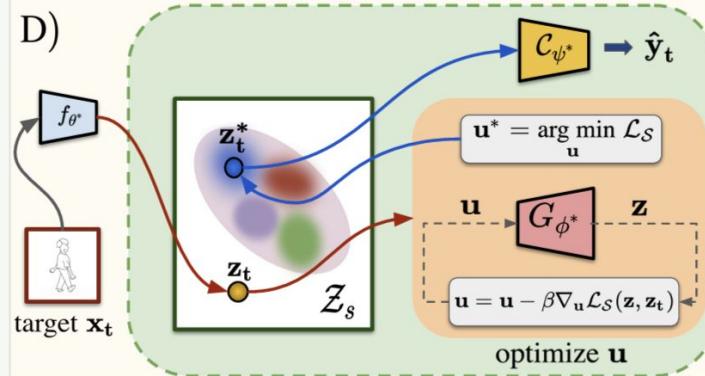
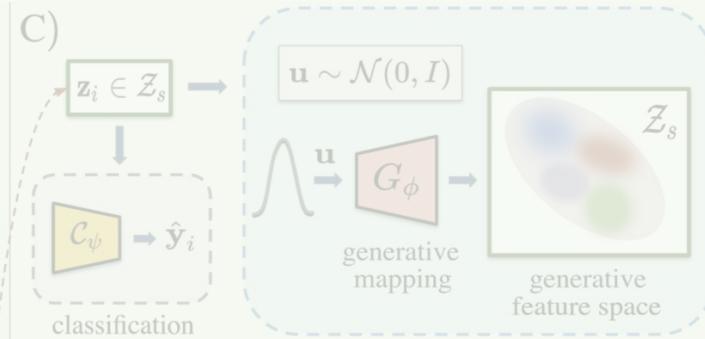
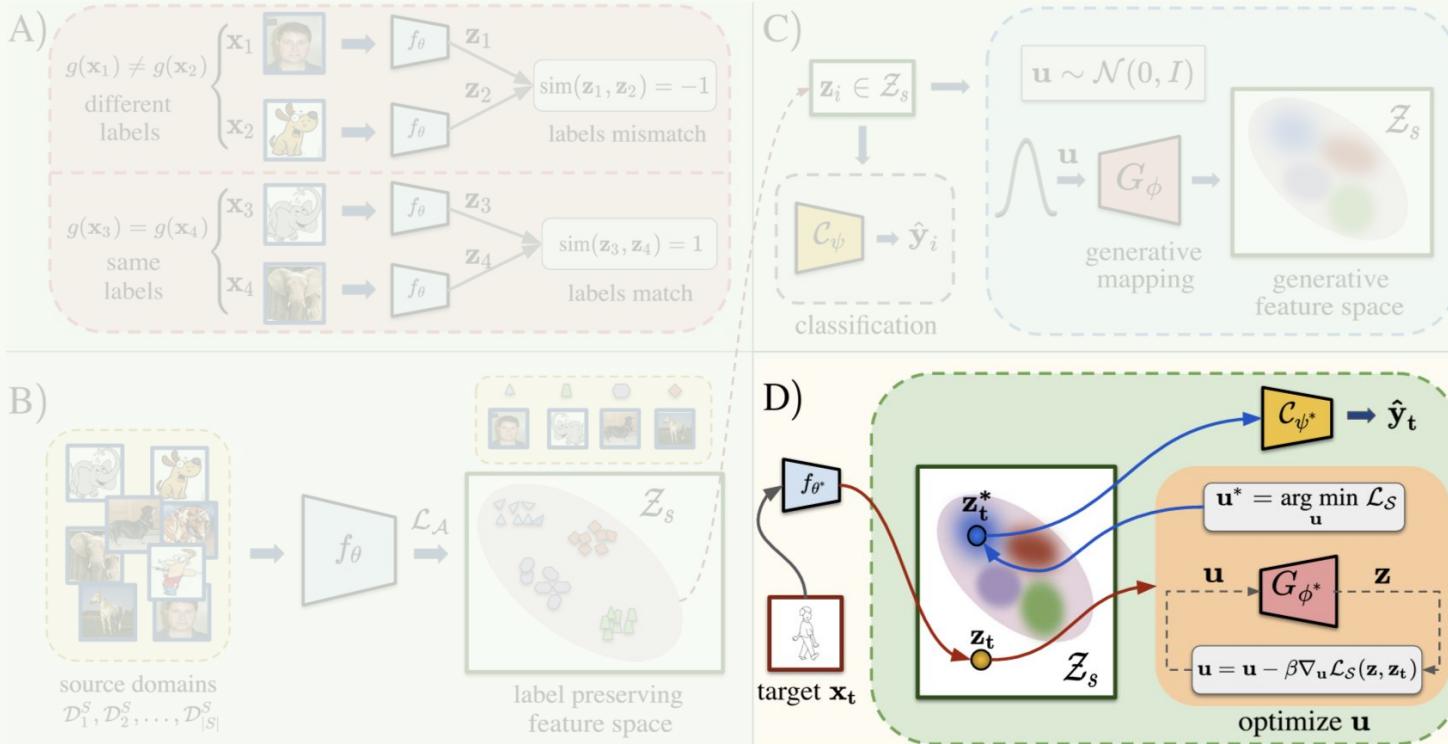
# Sample Adaptation



# Sample Adaptation



# Sample Adaptation



# Sample Adaptation

Method	MNIST	MNIST-M	SVHN	SYN	Avg.
Deep All	95.24±0.1	58.36±0.6	62.12±0.5	78.94±0.3	73.66
Jigen [5]	96.50±0.0	61.40±0.0	63.70±0.0	74.00±0.0	73.90
CCSA [29]	95.20±0.0	58.20±0.0	65.50±0.0	79.10±0.0	74.50
MMD-AAE [24]	96.50±0.0	58.40±0.0	65.00±0.0	78.40±0.0	74.60
CrossGrad [38]	96.70±0.0	61.10±0.0	65.30±0.0	80.20±0.0	75.80
L2A-OT [47]	96.70±0.0	63.90±0.0	68.60±0.0	83.20±0.0	78.10
Ours	<b>97.99±0.1</b>	<b>66.52±0.4</b>	<b>71.31±0.3</b>	<b>85.40±0.5</b>	<b>80.30</b>

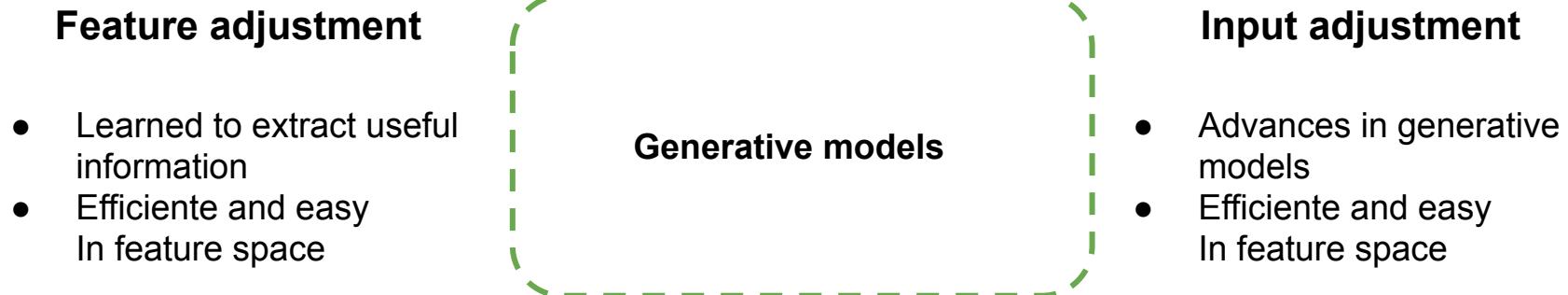
Method	Artistic	Clipart	Product	Real-World	Avg.
Deep All	52.06±0.5	46.12±0.3	70.45±0.2	72.45±0.2	60.27
D-SAM [7]	58.03±0.0	44.37±0.0	69.22±0.0	71.45±0.0	60.77
Jigen [5]	53.04±0.0	47.51±0.0	71.47±0.0	72.79±0.0	61.20
MMD-AAE [24]	56.50±0.0	47.30±0.0	72.10±0.0	74.80±0.0	62.70
DSON [37]	59.37±0.0	45.70±0.0	71.84±0.0	74.68±0.0	62.90
RSC [16]	58.42±0.0	47.90±0.0	71.63±0.0	74.54±0.0	63.12
L2A-OT [47]	60.60±0.0	50.10±0.0	74.80±0.0	77.00±0.0	65.60
Ours	<b>62.63±0.2</b>	<b>55.79±0.3</b>	<b>76.86±0.1</b>	<b>78.98±0.1</b>	<b>68.56</b>

Method	Caltech	LabelMe	Pascal	Sun	Avg.
Deep All	96.45±0.1	60.03±0.5	70.41±0.4	62.63±0.3	72.38
Jigen [5]	96.93±0.0	60.90±0.0	70.62±0.0	64.30±0.0	73.19
MMLD [27]	96.66±0.0	58.77±0.0	71.96±0.0	68.13±0.0	73.88
MASF [6]	94.78±0.2	64.90±0.1	69.14±0.2	67.64±0.1	74.11
EISNet [45]	97.33±0.4	63.49±0.8	69.83±0.5	68.02±0.8	74.67
RSC [16]	97.61±0.0	61.86±0.0	73.93±0.0	68.32±0.0	75.43
Ours	<b>98.12±0.1</b>	<b>66.80±0.3</b>	<b>74.77±0.4</b>	<b>70.43±0.1</b>	<b>77.53</b>

Method	Caltech	LabelMe	Pascal	Sun	Avg.
Deep All	96.45	60.03	70.41	62.63	72.38
1-NN	96.51	61.44	71.82	63.46	73.31
Ours ( $G_\phi$ = GAN)	97.89	<b>67.18</b>	74.59	70.28	77.48
Ours ( $G_\phi$ = VAE)	<b>98.12</b>	66.80	<b>74.77</b>	<b>70.43</b>	<b>77.53</b>

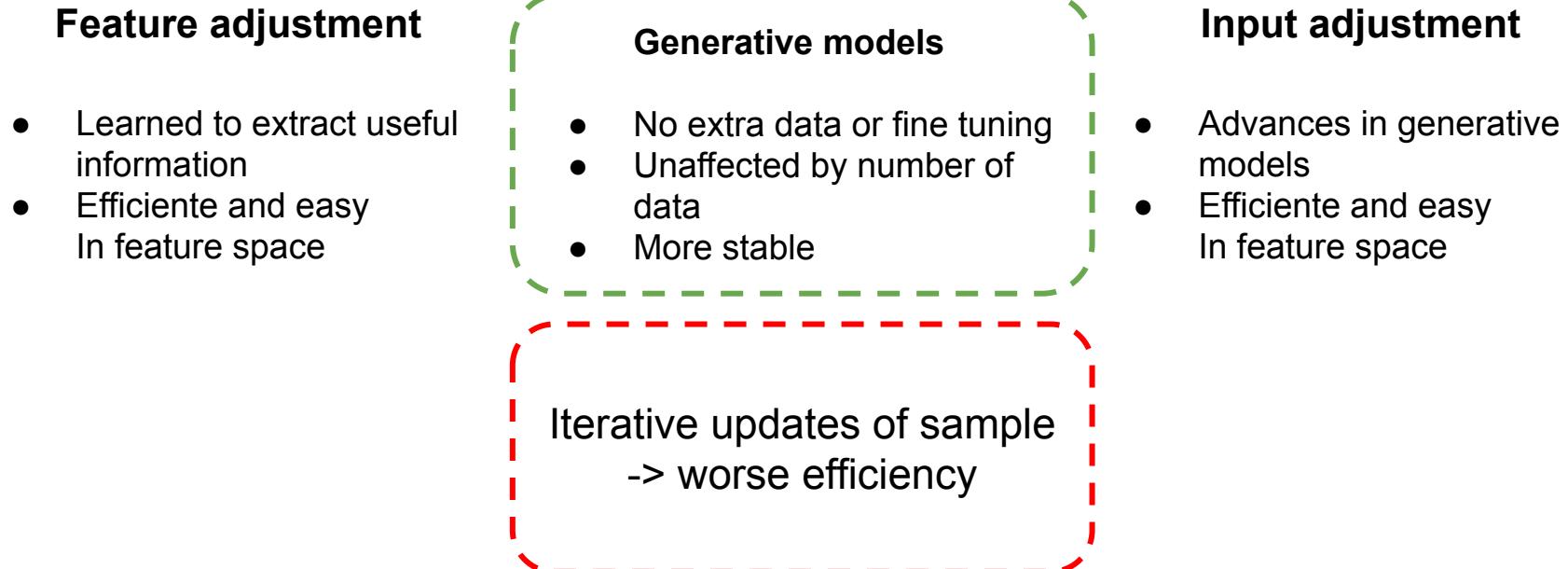
# Sample Adaptation

$$\mathbf{x}'_s = \psi(\mathbf{x}_t), \quad \mathbf{y}_t = f_{\theta_s}(\mathbf{x}'_s).$$

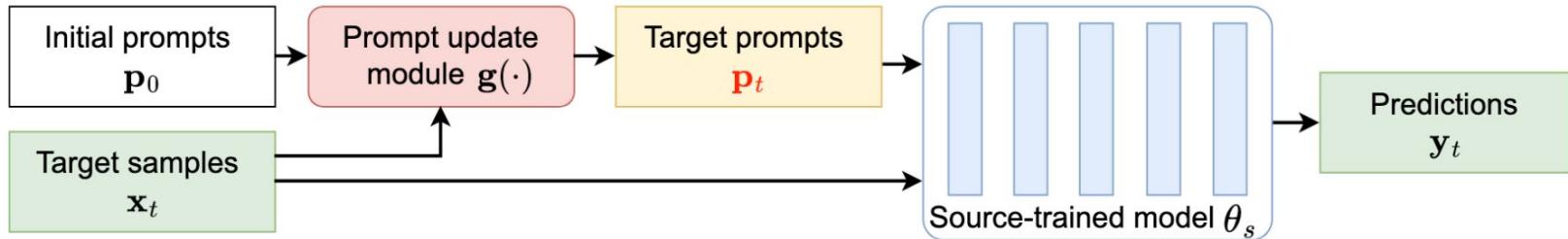


# Sample Adaptation

$$\mathbf{x}'_s = \psi(\mathbf{x}_t), \quad \mathbf{y}_t = f_{\theta_s}(\mathbf{x}'_s).$$



# Prompt Adaptation



$$\mathbf{p}_t = \mathbf{g}(\mathbf{x}_t, \mathbf{p}_0), \quad \mathbf{y}_t = f_{\theta}(\mathbf{x}_t, \mathbf{p}_t)$$

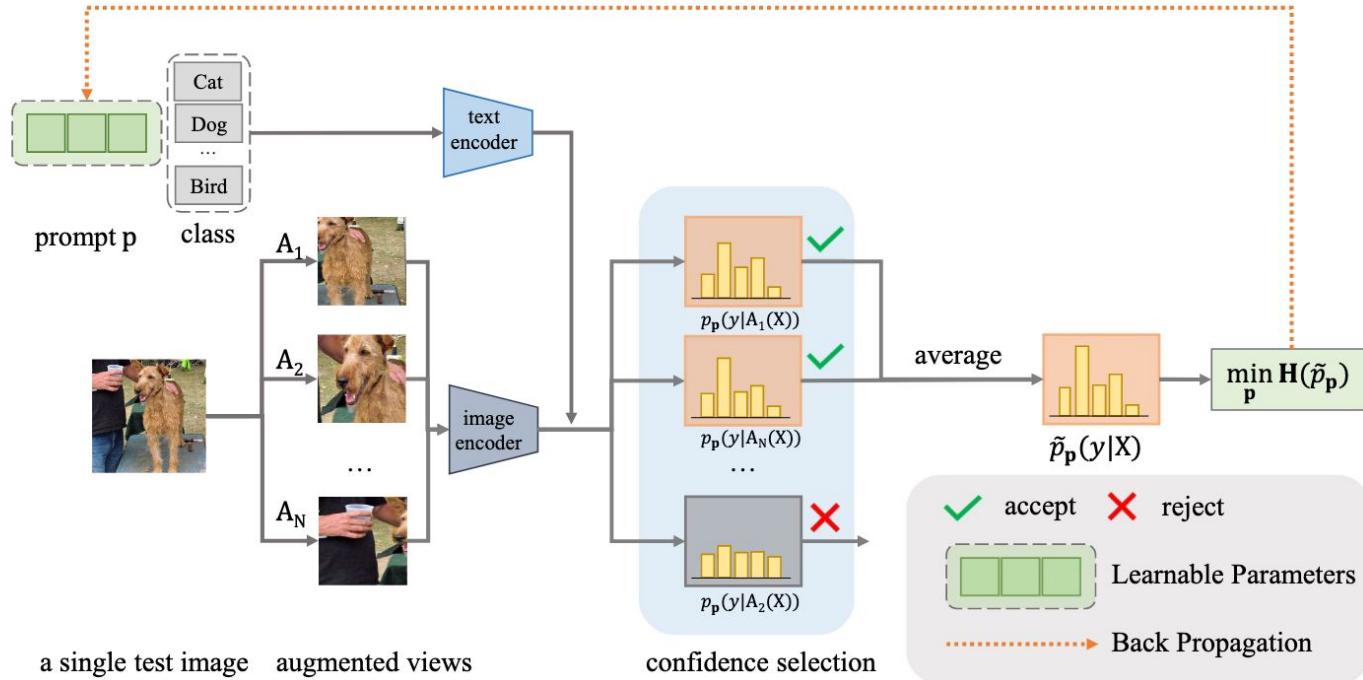
## Text-prompting

- Prompt-drive approaches
- Use LLMs to provide description or concepts

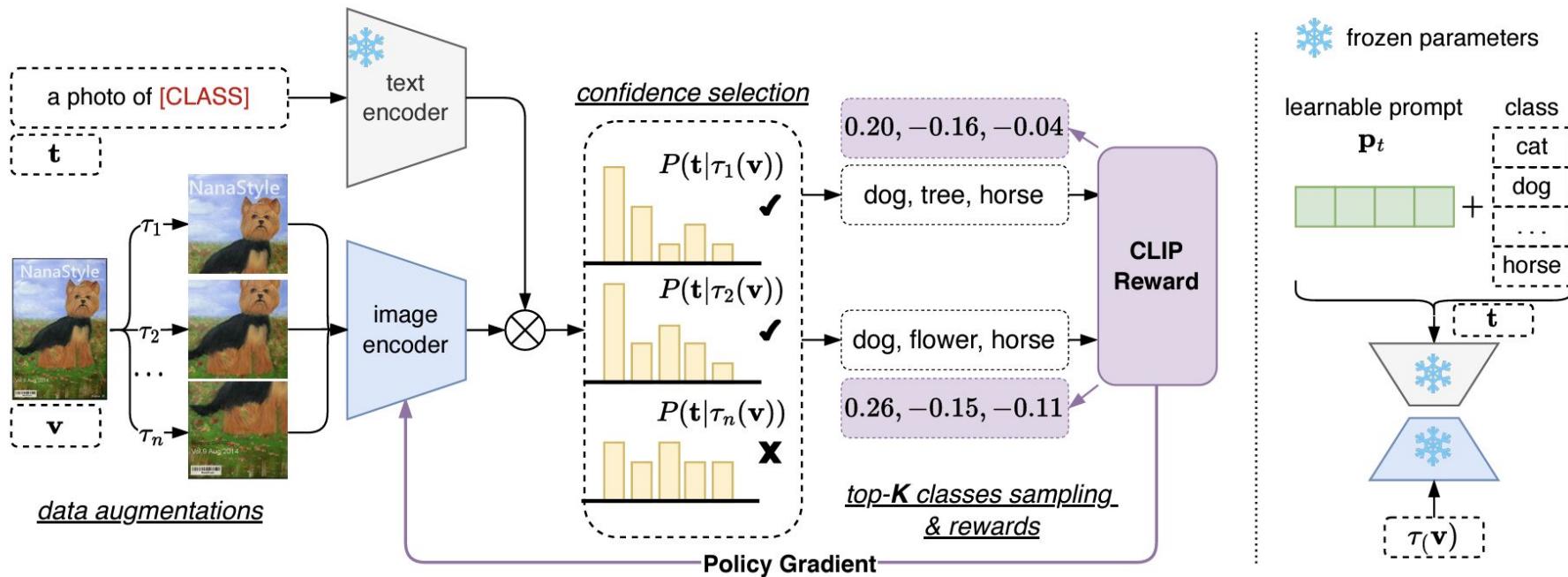
## Embedding-prompting

- Leverage Transformer and VLMs

# Prompt Adaptation



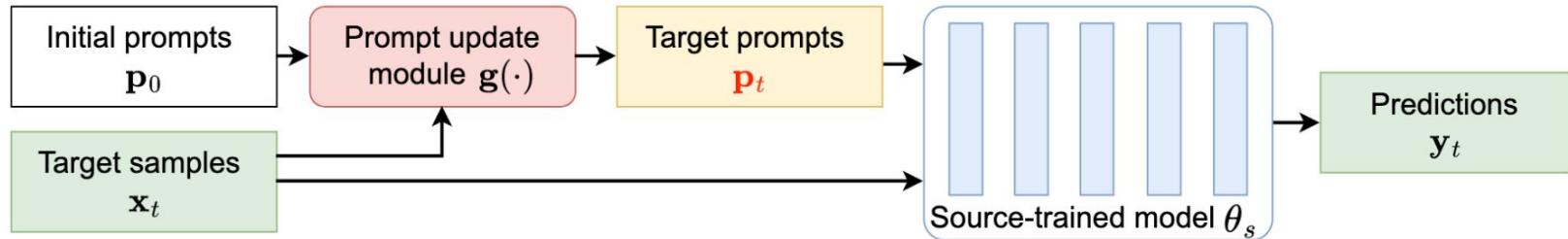
# Prompt Adaptation



# Prompt Adaptation

Method	ImageNet	ImageNet-A	ImageNet-V2	ImageNet-R	ImageNet-Sketch	OOD Average
<i>Zero-shot baseline</i>						
CLIP-ViT-B/16	66.73	47.87	60.86	73.98	46.09	57.20
CLIP-ViT-L/14	73.44	68.82	67.80	85.40	57.84	69.97
<i>Prompt tuning for CLIP-ViT-B/16</i>						
CoOp (Zhou et al., 2021)	71.51	49.71	64.20	75.21	47.99	59.28
CoCoOp (Zhou et al., 2022)	71.02	50.63	64.07	76.18	48.75	59.91
TPT (Manli et al., 2022)	68.98	54.77	63.45	77.06	47.94	60.81
TPT + CoOp (Manli et al.)	73.61	57.95	66.83	77.27	49.29	62.84
TPT + CoOp + KD (Hinton et al.)	71.40	63.25	65.28	82.70	55.78	66.75
<b>RLCF</b>	73.23(↑6.50)	65.45(↑17.58)	69.77(↑8.91)	83.35(↑9.37)	54.74(↑8.65)	68.33(↑11.13)
<b>RLCF + CoOp</b>	76.05(↑4.54)	69.74(↑20.03)	70.62(↑6.42)	84.51(↑9.30)	56.49(↑8.50)	70.34(↑11.06)
<b>RLCF-S + CoOp</b>	76.50(↑4.99)	71.11(↑21.40)	70.92(↑6.72)	84.73(↑9.52)	56.97(↑8.98)	70.93(↑11.65)
<i>Image encoder tuning for CLIP-ViT-B/16</i>						
Pseudo-label (Lee et al., 2013)	69.11	62.15	63.56	80.03	49.45	63.80
TPT (Manli et al., 2022)	69.42	61.62	63.70	79.74	49.47	63.63
KD (Hinton et al., 2015)	70.92	66.39	65.01	82.12	53.51	66.76
ATKD (Guo et al., 2020)	70.51	70.66	65.54	85.12	53.56	68.72
<b>RLCF</b>	74.85(↑8.12)	73.71(↑25.84)	69.77(↑8.91)	86.19(↑12.21)	57.10(↑11.01)	71.69(↑14.49)
<b>RLCF-S</b>	75.34(↑8.61)	75.00(↑27.13)	70.08(↑9.22)	86.97(↑12.99)	57.75(↑11.66)	72.45(↑15.25)
<b>RLCF-S-M</b>	75.48(↑8.75)	75.16(↑27.29)	70.42(↑9.56)	87.23(↑13.25)	57.73(↑11.64)	72.64(↑15.44)

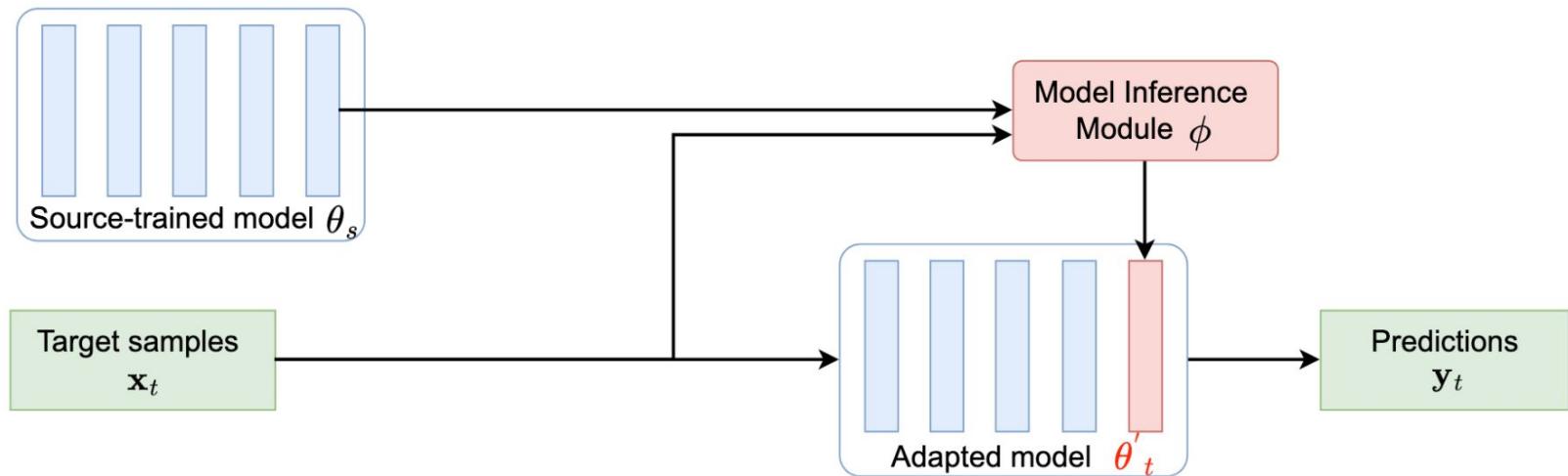
# Prompt Adaptation



$$\mathbf{p}_t = \mathbf{g}(\mathbf{x}_t, \mathbf{p}_0), \quad \mathbf{y}_t = f_{\theta}(\mathbf{x}_t, \mathbf{p}_t)$$

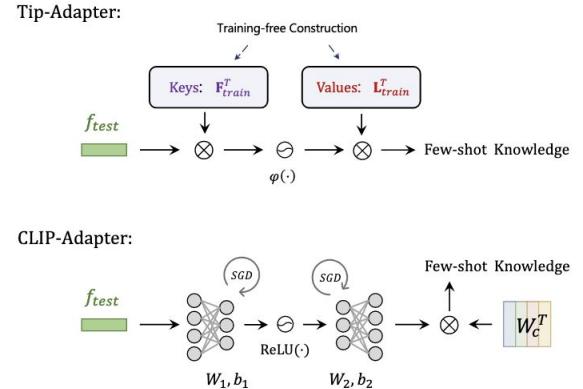
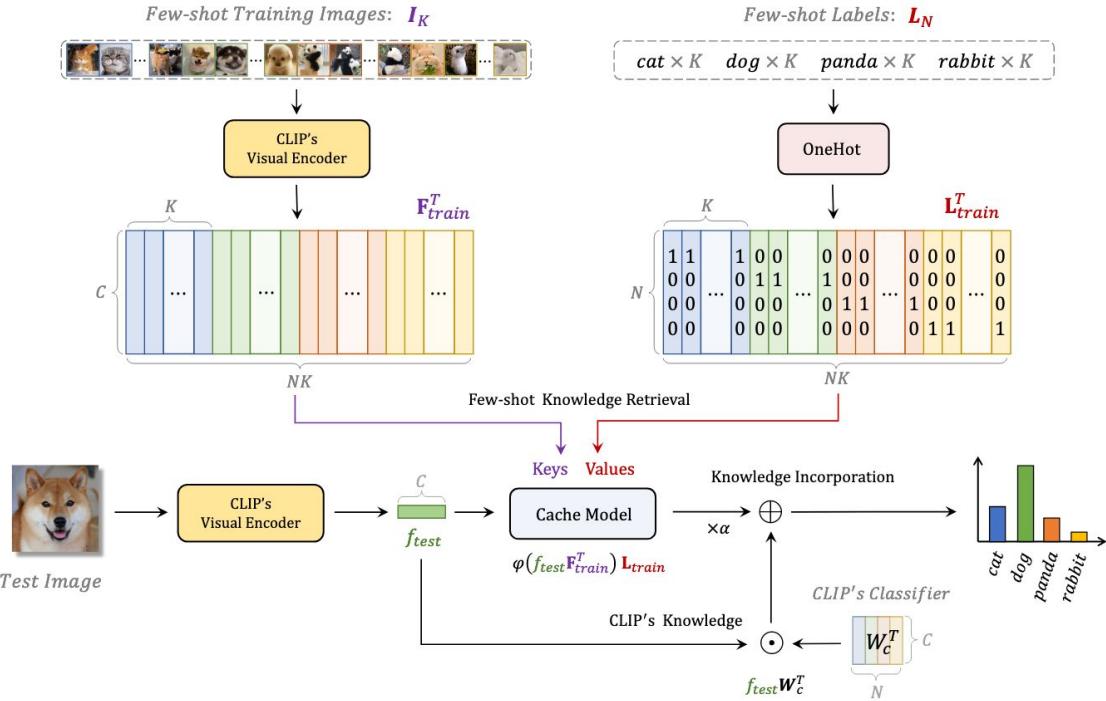
- may be overconfident
- NLP drawbacks
- task-specific

# Inference Adaptation

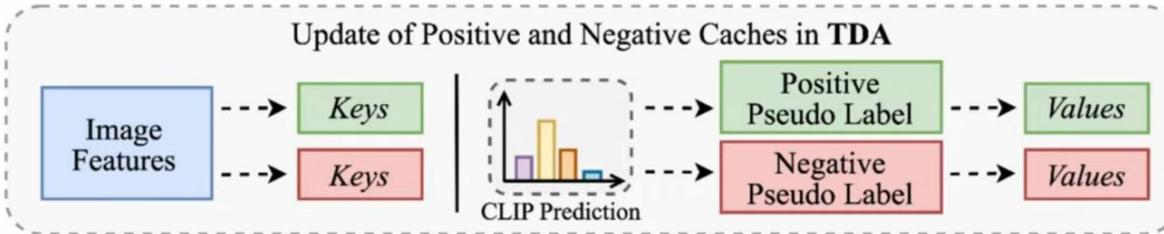


$$\theta'_t = \phi(\mathbf{x}_t, \theta_s), \quad \mathbf{y}_t = f_{\theta'_t}(\mathbf{x}_t)$$

# Tip-Adapter



# Training-free Dynamic Adapter (TDA)



- Cache Update Mechanism:

Both key-value caches update dynamically, integrating test features and corresponding pseudo labels from CLIP predictions, based on prediction entropy and cache capacity

## Positive Cache

- Capacity:**  
Holds up to  $k$  high-confidence samples per class
- Impact:**  
Continuously improves testing accuracy by integrating the most reliable samples

## Negative Cache

- Capacity:**  
Contains up to  $k$  samples per class, focused on uncertain predictions
- Impact:**  
Actively reduces prediction errors by managing and adjusting uncertain pseudo labels

## Proposal 1: Create a neutral cache

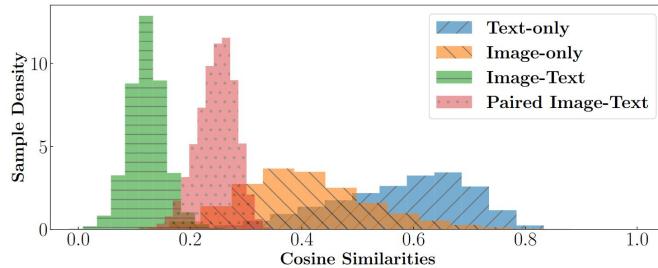
- Idea: Intermediate cache between positive and negative caches.

$$P_{\text{TDA}}(f_{\text{test}}) = f_{\text{test}} \mathbf{W}_c^T + P_{\text{pos}}(f_{\text{test}}) + P_{\text{neg}}(f_{\text{test}})$$

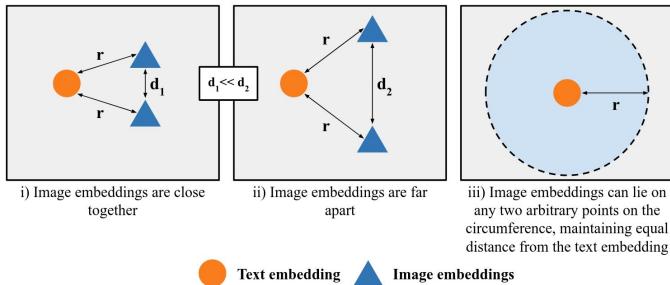


$$P_{\text{TDA}}(f_{\text{test}}) = f_{\text{test}} W_c^T + P_{\text{pos}}(f_{\text{test}}) + P_{\text{neg}}(f_{\text{test}}) + P_{\text{neut}}(f_{\text{test}})$$

# Proposal 2: Tackle intra-modality gap with TIP-X



**(a) Intra-modal and inter-modal CLIP cosine similarities.** We observe quite distinct intra-modal and inter-modal cosine similarity distributions.



**(b) Intra-modal degrees of freedom.** Different intra-modal similarities can satisfy same inter-modal constraints, leaving room for poor calibration.

Idea: Compute similarities between image features and text to decrease intra-modality gap.

What about using support sets obtained from LAION-5B?

# Proposal 2: Tackle intra-modality gap with TIP-X

$C$ : number of classes

$K$ : number of examples of classes in our cache

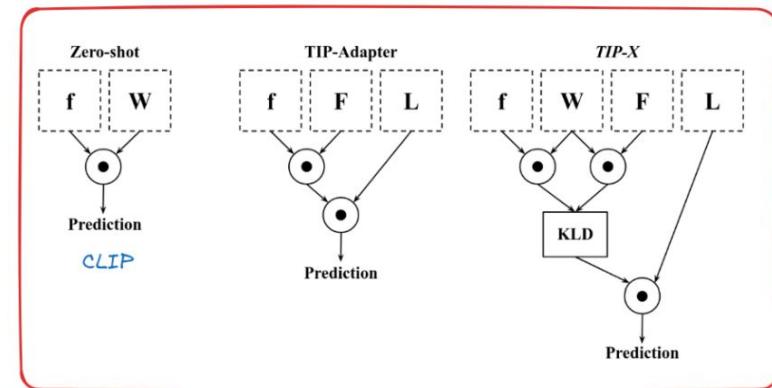
$$S = \text{softmax}(FW^T), S \in \mathbb{R}^{CK \times C}$$

$$s = \text{softmax}(fW^T), s \in \mathbb{R}^{t \times C}$$

$$M_{i,j} = \text{KL}(s_i || S_j), i \in [1, t], j \in [1, CK]$$

$$\text{KL}(P || Q) = \sum_i P_i \log \frac{P_i}{Q_i}$$

discrete probabilities  $\rightarrow$  discrete KL-div



$$\text{TXL} = fW^T + \alpha AL + \gamma \psi(-M)L$$

TIP-X Logits

$\psi(-M)L$

Rescales  $M$  to have same min max values as the TIP-Adapter affinities

inverts similarity meaning

Thank you for your attention

Questions?