



# Black-Box Forgetting

**Yusuke Kuwana<sup>1</sup> Yuta Goto<sup>1</sup> Takashi Shibata<sup>2</sup> Go Irie<sup>1</sup>**

**<sup>1</sup>Tokyo University of Science    <sup>2</sup>NEC Corporation**

# Motivation

Large-scale pre-trained models (PTMs) have strong capabilities of zero-shot classification for everyday objects.

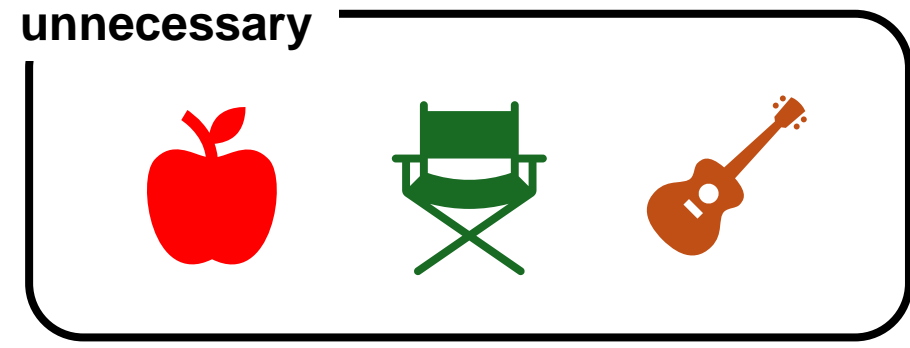
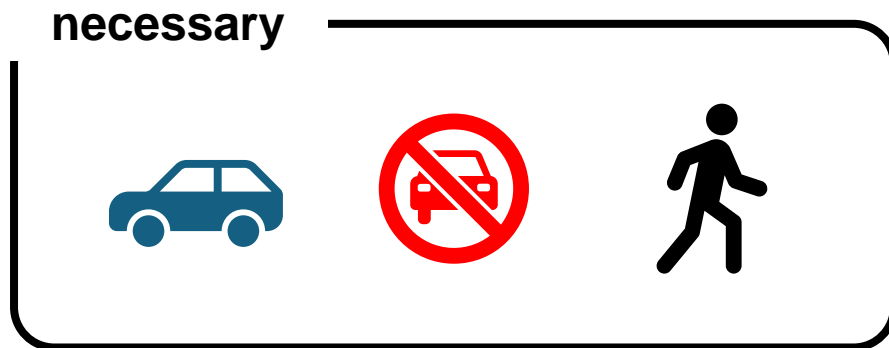
ex. CLIP [Radford+, ICML21]

	Dataset Examples	ImageNet ResNet101	Zero-Shot CLIP	$\Delta$ Score
ImageNet		76.2	76.2	0%
ObjectNet		32.6	72.3	+39.7%
ImageNet Sketch		25.2	60.2	+35.0%

# Motivation

Practical applications do not always require the classification of all kinds of objects.

ex. Autonomous driving system



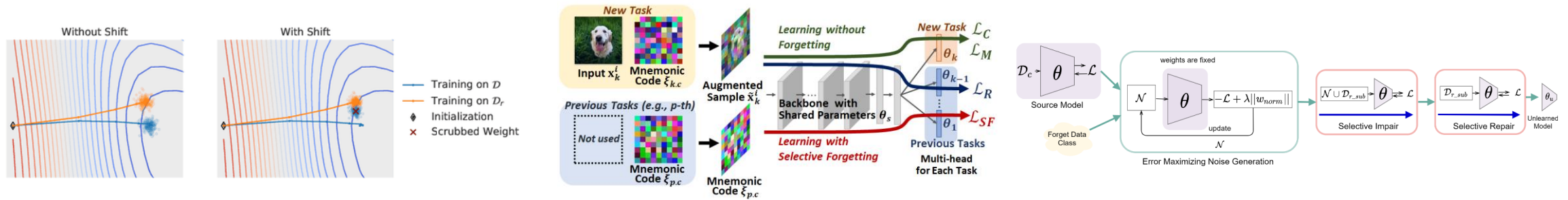
Retaining unnecessary classes may have disadvantages.

- Decrease overall accuracy
- Information leakage

**We address the problem of selective forgetting.**

# Motivation

## Existing selective forgetting methods



[Golatker+, ECCV20]

[Shibata+, IJCAI21]

[Tarun+, IEEE TNNLS23]

All the existing methods assume “white-box” settings, where model information is available for training.

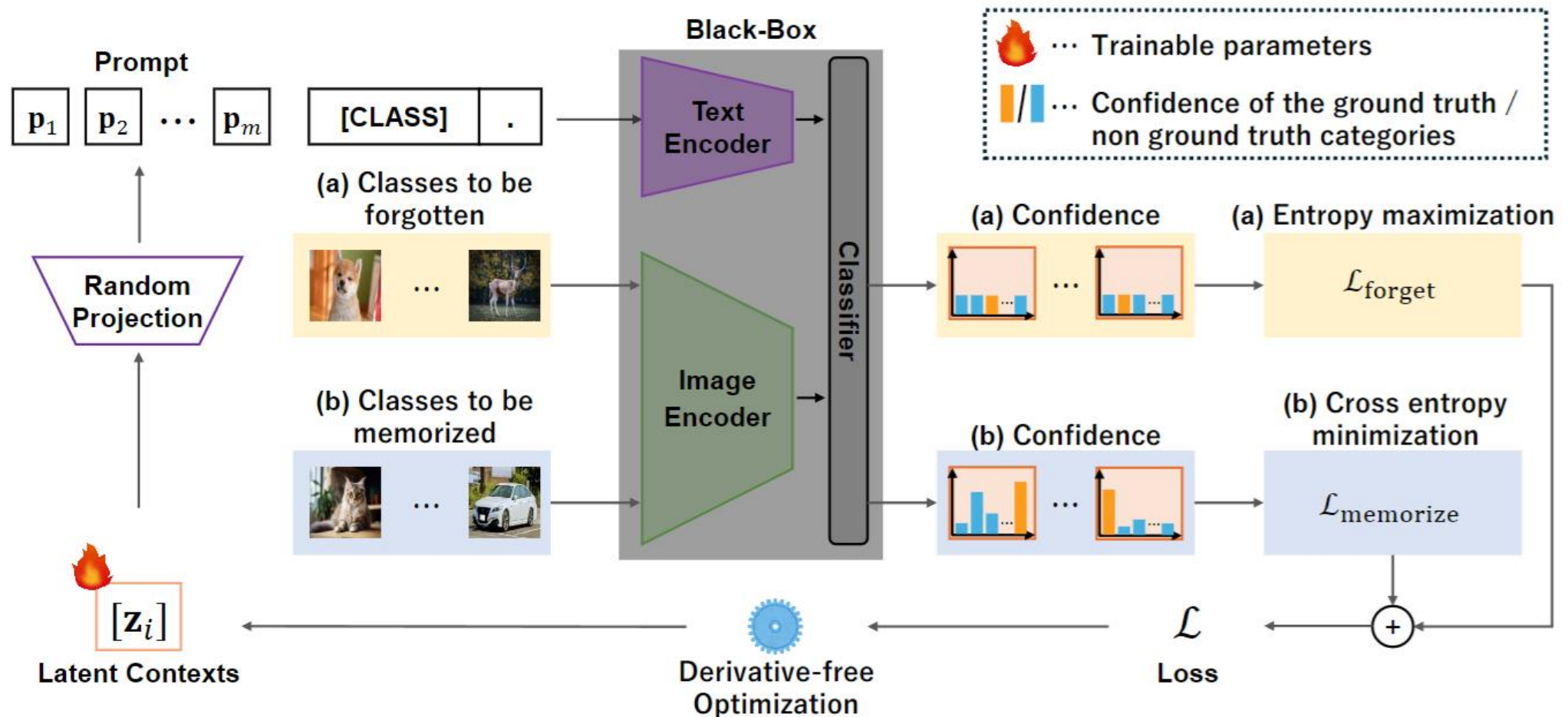
However, PTMs are often “black-box,” where model information is unavailable for commercial reasons or social responsibilities.

→ **Inapplicable for black-box models.**

**We address a novel problem of selective forgetting for black-box models.**

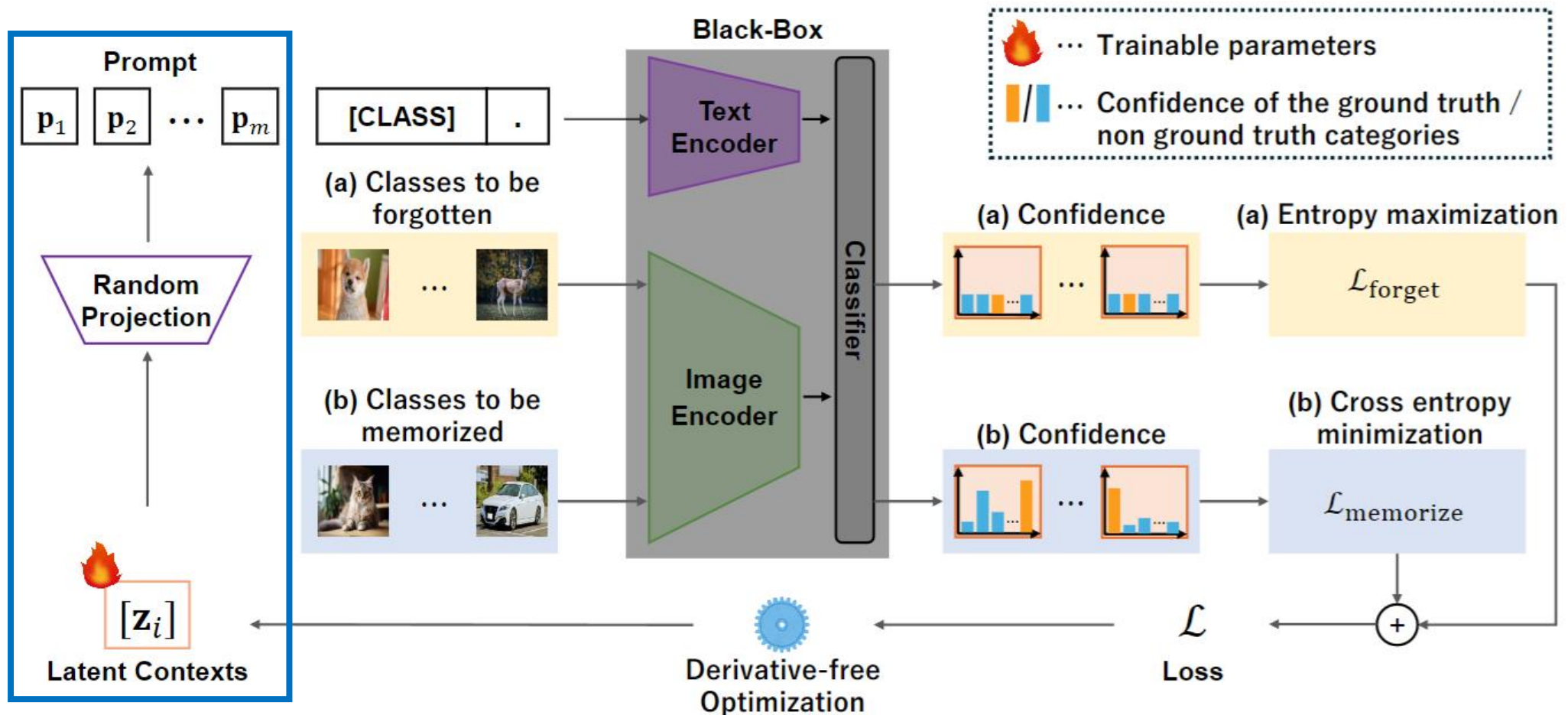
# Our Method

We optimize the textual prompt to decrease the accuracy of specified classes through derivative-free optimization, because the gradients of the loss are unavailable in black-box models.



# Our Method

We optimize lower-dimensional latent contexts instead of optimizing contexts for the textual prompt directly to mitigate high-dimensional optimization.



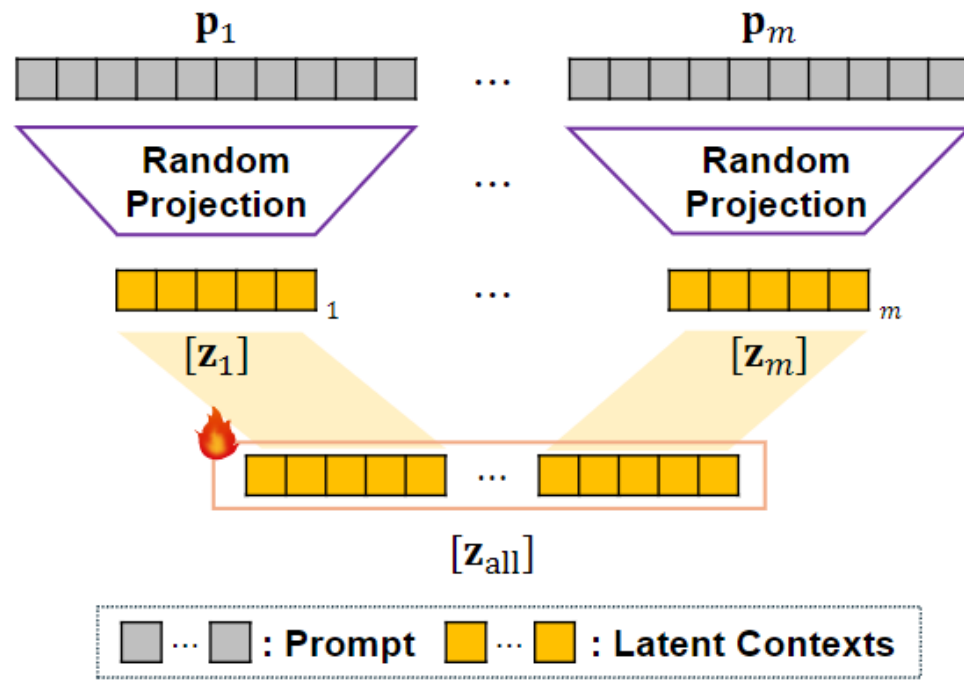


# Our Method

BBT [Sun et al., ICML22] optimizes a lower-dimensional latent context instead of directly optimizing textual prompt to mitigate high dimensionality.

We found that the effectiveness of context parametrization through BBT is **limited** for Black-Box Forgetting.

BBT [Sun et al., ICML22]



Method	CIFAR-10		
	$H \uparrow$	$Err_{\text{for}} \uparrow$	$Acc_{\text{mem}} \uparrow$
Zero-Shot CLIP	15.30	8.37	89.05
BBT [Sun et al., ICML22]	<b>85.69</b>	<b>79.31</b>	<b>93.19</b>
CoOp [Zhou et al., IJCV22] (White-Box)	96.49	96.95	96.04

$Err_{\text{for}}$  ... the error of the classes to be forgotten

$Acc_{\text{mem}}$  ... the accuracy of the classes to be memorized

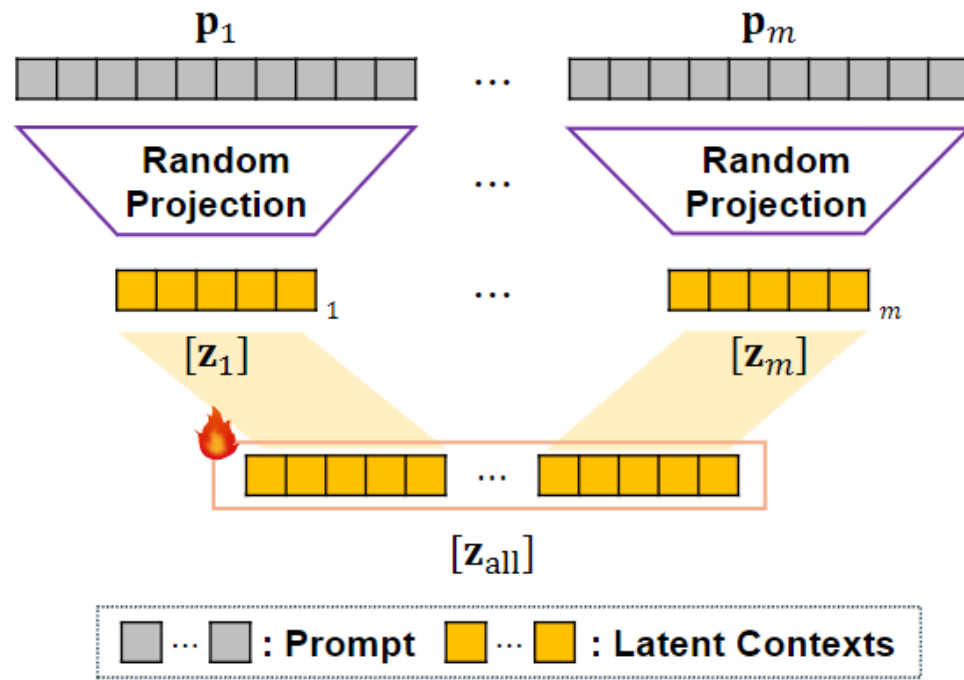
$H$  ... the harmonic mean of  $Err_{\text{for}}$  and  $Acc_{\text{mem}}$

# Our Method

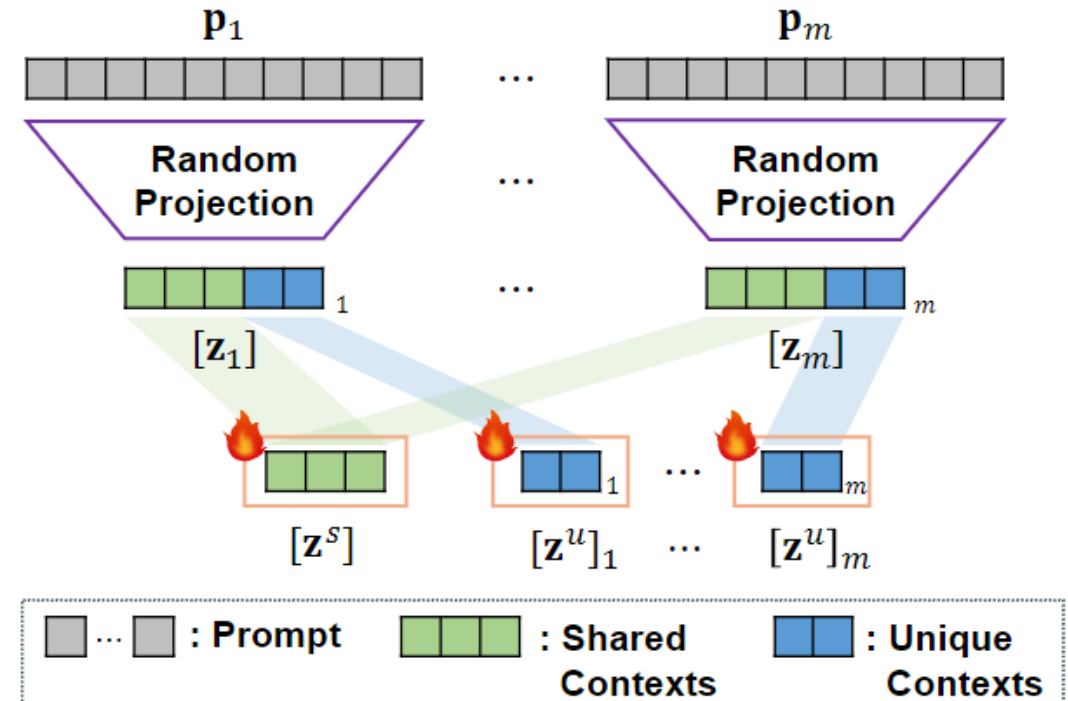
We propose **Latent Context Sharing (LCS)** for more effective context parametrization.

In **LCS**, a latent context is composed of **unique** components and **common** components among multiple latent contexts, and each component is optimized independently.

BBT [Sun et al., ICML22]



**Ours: Latent Context Sharing (LCS)**





# Experiments

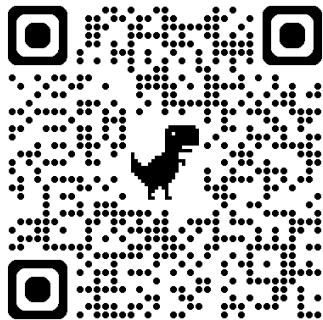
Our method outperforms the reasonable baselines on all the datasets.

**Table 1: Comparisons with the baselines.** The best value is shown in bold. BBT [Sun et al., 2022b] and CBBT (w/o adapter) [Guo et al., 2023] are the reasonable baselines as these are for black-box prompt tuning. CoOp [Zhou et al., 2022b] is a white-box method and is included for a reference. Performance is evaluated using the three metrics: the error  $Err_{\text{for}}$  for the classes to be forgotten, the accuracy  $Acc_{\text{mem}}$  for the classes to be memorized, the harmonic mean  $H$  of  $Err_{\text{for}}$  and  $Acc_{\text{mem}}$ . Higher values mean better performance.

Method	CIFAR-10			CIFAR-100		
	$H \uparrow$	$Err_{\text{for}} \uparrow$	$Acc_{\text{mem}} \uparrow$	$H \uparrow$	$Err_{\text{for}} \uparrow$	$Acc_{\text{mem}} \uparrow$
Zero-Shot CLIP	15.30	8.37	89.05	42.14	31.17	65.03
BBT	85.69 $\pm$ 0.02	79.31 $\pm$ 0.03	93.19 $\pm$ 0.01	78.36 $\pm$ 0.01	87.30 $\pm$ 0.01	71.09 $\pm$ 0.00
CBBT	93.48 $\pm$ 0.02	90.99 $\pm$ 0.04	<b>96.11</b> $\pm$ 0.00	73.20 $\pm$ 0.00	72.69 $\pm$ 0.01	<b>73.72</b> $\pm$ 0.00
Ours (w/o LCS)	72.37 $\pm$ 0.13	58.57 $\pm$ 0.17	94.68 $\pm$ 0.01	79.38 $\pm$ 0.02	89.17 $\pm$ 0.03	71.52 $\pm$ 0.01
Ours	<b>95.07</b> $\pm$ 0.01	<b>96.10</b> $\pm$ 0.02	94.06 $\pm$ 0.01	<b>80.99</b> $\pm$ 0.01	<b>93.37</b> $\pm$ 0.02	71.52 $\pm$ 0.01
CoOp (White-Box)	96.49 $\pm$ 0.00	96.95 $\pm$ 0.01	96.04 $\pm$ 0.00	82.22 $\pm$ 0.00	99.81 $\pm$ 0.00	69.90 $\pm$ 0.01
Method	CUB-200-2011			ImageNet30		
	$H \uparrow$	$Err_{\text{for}} \uparrow$	$Acc_{\text{mem}} \uparrow$	$H \uparrow$	$Err_{\text{for}} \uparrow$	$Acc_{\text{mem}} \uparrow$
Zero-Shot CLIP	46.30	46.20	<b>46.41</b>	2.31	1.17	98.00
BBT	58.75 $\pm$ 0.01	88.98 $\pm$ 0.04	43.85 $\pm$ 0.01	94.22 $\pm$ 0.05	90.17 $\pm$ 0.08	99.06 $\pm$ 0.01
CBBT	56.84 $\pm$ 0.01	73.52 $\pm$ 0.02	46.33 $\pm$ 0.01	87.88 $\pm$ 0.08	79.69 $\pm$ 0.12	<b>99.32</b> $\pm$ 0.02
Ours (w/o LCS)	58.78 $\pm$ 0.01	85.85 $\pm$ 0.01	44.69 $\pm$ 0.01	95.26 $\pm$ 0.02	92.19 $\pm$ 0.03	98.59 $\pm$ 0.01
Ours	<b>59.67</b> $\pm$ 0.01	<b>89.29</b> $\pm$ 0.01	44.81 $\pm$ 0.01	<b>97.28</b> $\pm$ 0.01	<b>95.94</b> $\pm$ 0.01	98.67 $\pm$ 0.01
CoOp (White-Box)	63.20 $\pm$ 0.02	98.09 $\pm$ 0.02	46.62 $\pm$ 0.02	99.30 $\pm$ 0.01	99.72 $\pm$ 0.00	98.89 $\pm$ 0.01

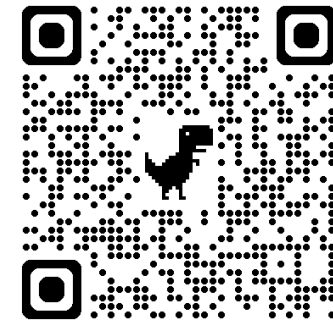
# Conclusion

- We proposed **Black-Box Forgetting**, a novel problem of selective forgetting for black-box models.
- We introduced **Latent Context Sharing (LCS)**, an efficient and effective parametrization method of prompt, which is suitable for derivative-free optimization.
- Experimental results demonstrated that our method outperforms the reasonable baselines.



arXiv

## Thank you!



code