

# InfoSAM: Fine-Tuning the Segment Anything Model from An Information- Theoretic Perspective

Yuanhong Zhang<sup>\*1</sup>, Muyao Yuan<sup>\*1</sup>, Weizhan Zhang<sup>1</sup>, Tieliang Gong<sup>1</sup>, Wen Wen<sup>1</sup>,  
Jiangyong Ying<sup>2</sup>, Weijie Shi<sup>1</sup>

<sup>1</sup> Xi'an Jiaotong University

<sup>2</sup> China Telecom E-surfing Vision Technology

<sup>\*</sup> Equal contribution

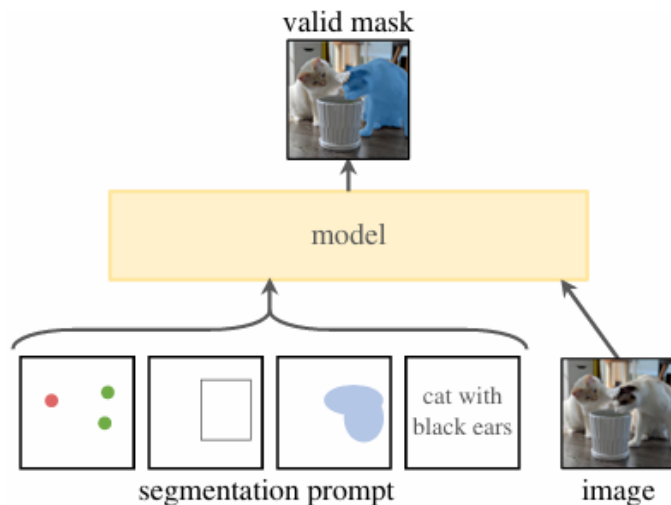
Project Page



# Background & Motivation

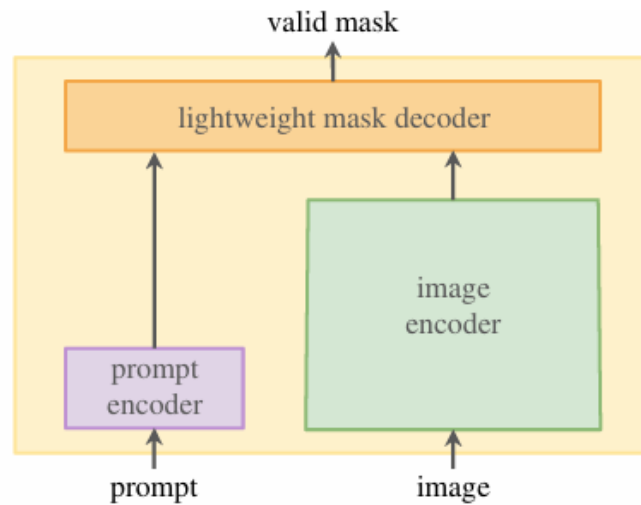


❑ **Segment Anything**<sup>[1]</sup>: shows impressive zero-shot performance on generic object segmentation



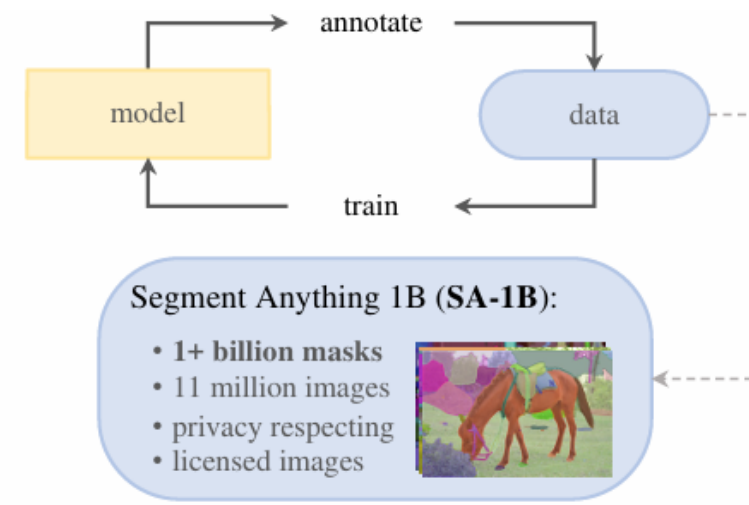
(a) **Task:** promptable segmentation

"interactive"



(b) **Model:** Segment Anything Model (SAM)

"Encoder-Decoder"



(c) **Data:** data engine (top) & dataset (bottom)

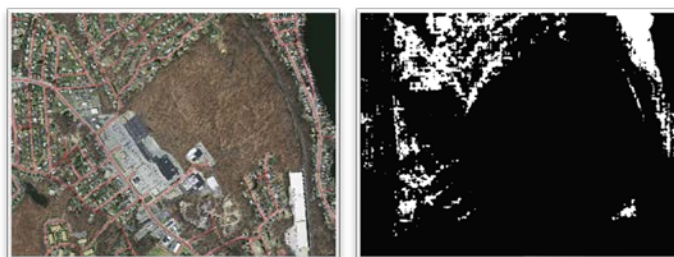
"over 1 billion masks"

❑ **Segment Anything:** still struggles with domain-specific real-world segmentation tasks

Camouflaged Scenes



Remote Sensing



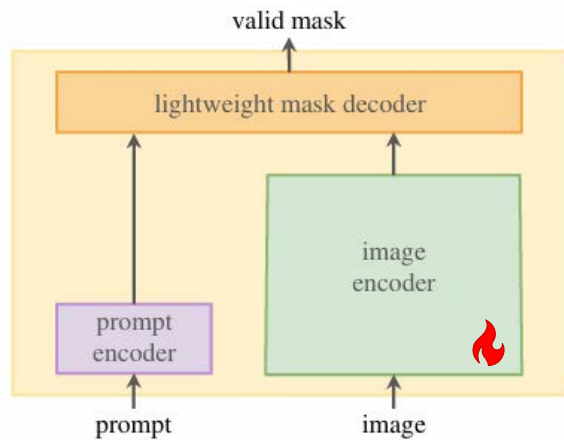
Agriculture



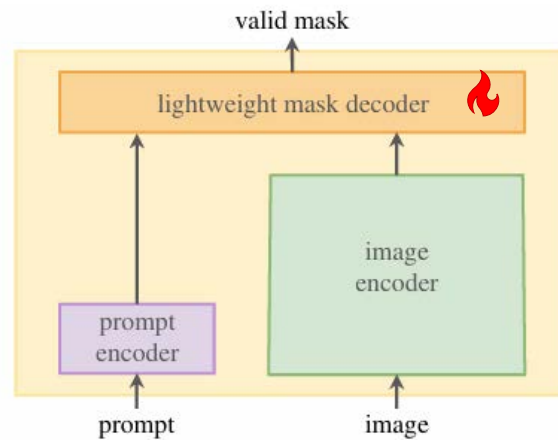
# Background & Motivation



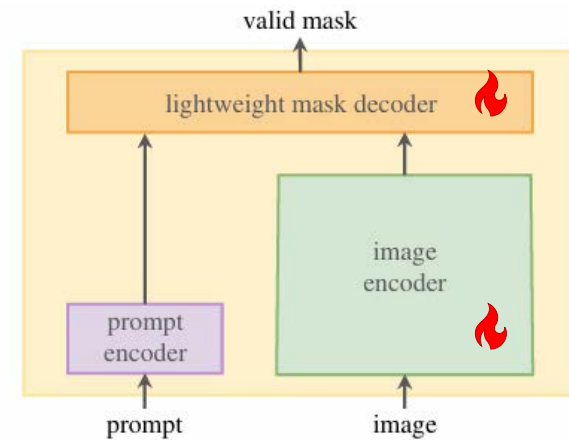
- ❑ **Parameter-Efficient Fine-Tuning (PEFT)**: a promising approach to unleash the potential of SAM in novel scenarios.



SAM-COBOT<sup>[2]</sup>



HQSAM<sup>[3]</sup>



CAT-SAM<sup>[4]</sup>

 Tunable

**"Ignore the beneficial information encoded in the pre-trained SAM !"**

- ❑ We argue that:
  - There exists **domain-invariant information** that emerges from extensive pre-training.
  - This information is embedded in the feature distributions **between the encoder and decoder**, yet it can be easily **overridden or suppressed** during fine-tuning.

# Challenges & Our solutions



## □ Challenges:

- How to **extract** a good domain-invariant information? → To compress
- How to effectively **transfer it** to the fine-tuned models? → To distill

## □ Formulation:

- Let  $z_i^T / z_i^S$  and  $z_m^T / z_m^S$  denote the image encoder features and mask decoder tokens from the teacher (pre-trained SAM) and student (fine-tuned SAM), respectively. We formulate the extraction and transfer process as the information flow. We use **matrix-based Rényi's theory** to quantify such information.
- **Objective 1:** to prioritize domain-invariant relations, we **constrain the information flow** via an upper bound  $I_c$ :

$$I_\alpha(z_i^T, z_m^T; r^T) \leq I_c$$

- **Objective 2:** To **maximize the extracted information** between pre-trained SAM and fine-tuned SAM:

$$\max_{\omega} I_\alpha(r^T; r^S)$$

- The Lagrangian formulation explicitly implements this trade-off:

$$\max_{\omega} I_\alpha(r^T; r^S) - \beta I_\alpha(z_i^T, z_m^T; r^T)$$

To distill

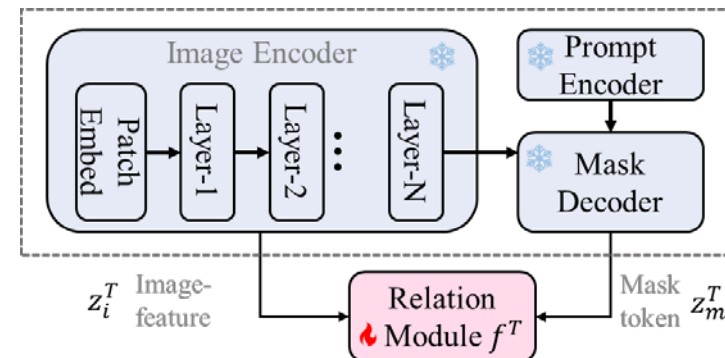
To compress

# Our solutions



## □ An Information View of SAM Distillation:

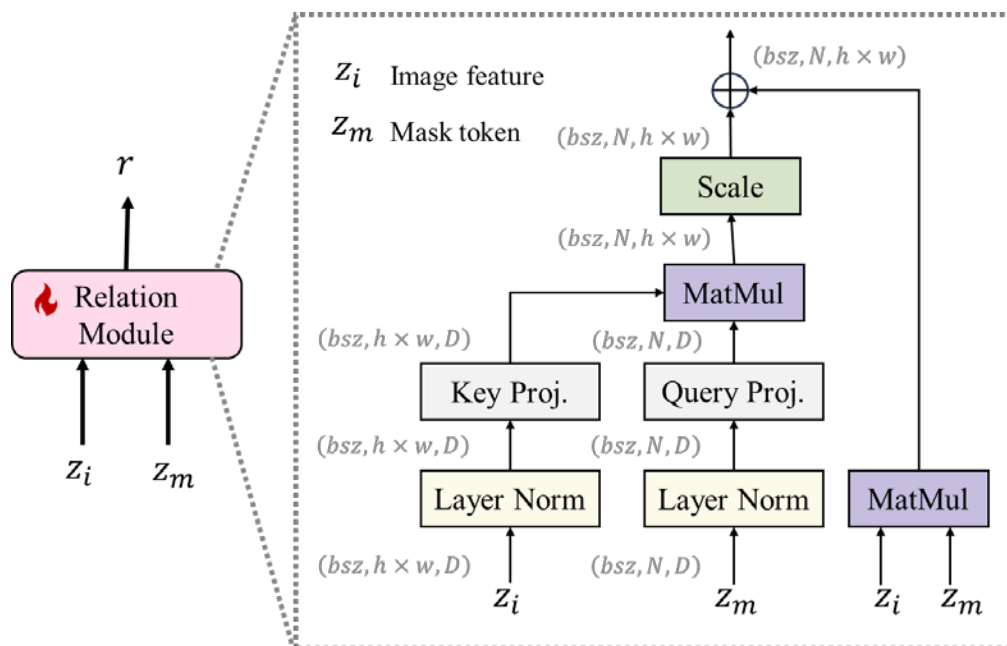
For Challenge 1: Compressing Intra-SAM Relations



Extraction & Compression

✓ Attention-based module  $f^T$  designed for extraction

✓ Guide relations toward domain-invariant cues



$$I_{\alpha}(z_i^T, z_m^T; r^T) \leq I_c$$



$$\begin{aligned} \mathcal{L}_r &= I_{\alpha}(z_i^T, z_m^T; r^T) \\ &= S_{\alpha}(G_i^T, G_m^T) + S_{\alpha}(G_r^T) - S_{\alpha}(G_i^T, G_m^T, G_r^T) \end{aligned}$$



Set  $\alpha = 2$

$$\mathcal{L}_r = -\log_2 \|G_r^T\|_F^2 + \log_2 \|G_{imr}^T\|_F^2$$

# Our solutions



## □ An Information View of SAM Distillation:

For Challenge 2: Maximizing Inter-SAM Relations

✓ Transfer the relationships by minimizing their distance.

$$\max_{\omega} I_{\alpha}(r^T; r^S)$$

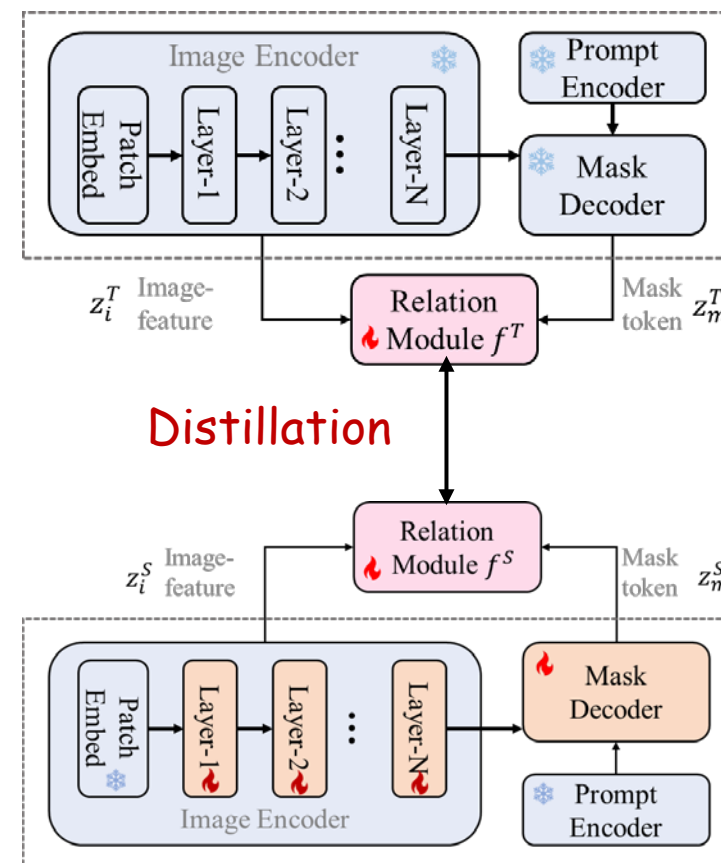


$$\begin{aligned}\mathcal{L}_d &= -I_{\alpha}(r^T; r^S) \\ &= -S_{\alpha}(G_r^T) - S_{\alpha}(G_r^S) + S_{\alpha}(G_r^T, G_r^S)\end{aligned}$$



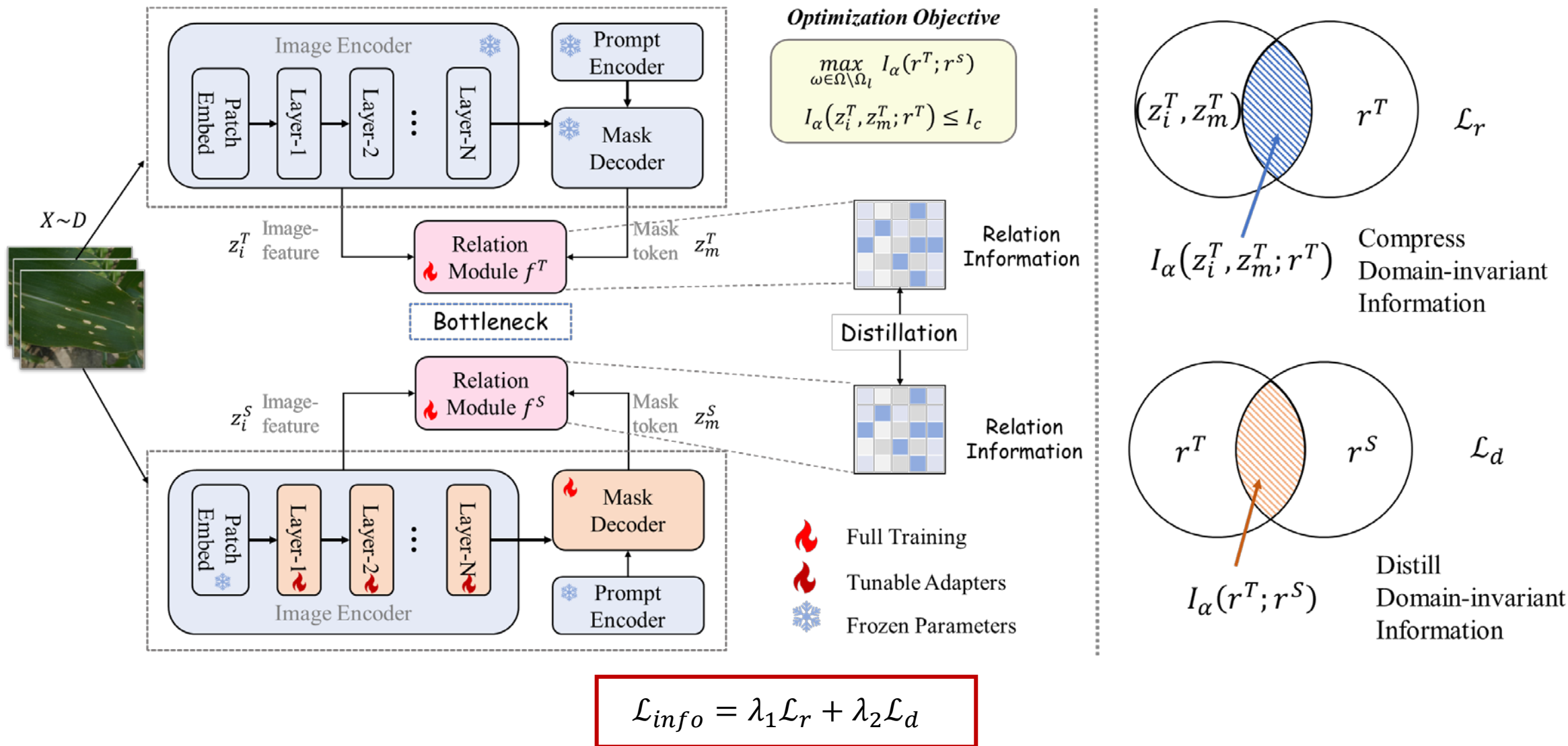
Set  $\alpha = 2$

$$\mathcal{L}_d = \log_2 \|G_r^T\|_F^2 + \log_2 \|G_r^S\|_F^2 - \log_2 \|G_r^{TS}\|_F^2$$



# Our solutions

## □ Overview of InfoSAM:





# Experiments



□ Compare with PEFT baselines across various downstream segmentation tasks

METHOD	NATURAL IMAGES			MEDICAL				AGRICULTURE		REMOTE SENSING	
	CAMO			ISIC 2017		Kvasir		Leaf		Road	
	$S_\alpha \uparrow$	$E_\phi \uparrow$	$F_\beta^\omega \uparrow$	Jac $\uparrow$	Dice $\uparrow$	$S_\alpha \uparrow$	$E_\phi \uparrow$	IoU $\uparrow$	Dice $\uparrow$	IoU $\uparrow$	Dice $\uparrow$
SAM decoder-only	$79.7 \pm 0.02$	$88.8 \pm 0.09$	$79.6 \pm 0.01$	$61.0 \pm 0.12$	$71.7 \pm 0.14$	$71.4 \pm 0.16$	$77.9 \pm 0.17$	$37.6 \pm 0.11$	$47.0 \pm 0.16$	$7.2 \pm 0.24$	$12.9 \pm 0.29$
	$84.9 \pm 0.38$	$92.7 \pm 0.34$	$81.8 \pm 0.33$	$85.9 \pm 0.34$	$92.2 \pm 0.20$	$90.9 \pm 0.05$	$95.2 \pm 0.18$	$55.6 \pm 1.12$	$68.8 \pm 1.17$	$47.6 \pm 0.47$	$64.1 \pm 0.47$
BitFit	$87.5 \pm 0.13$	$94.5 \pm 0.08$	$85.3 \pm 0.48$	$87.7 \pm 0.14$	$93.2 \pm 0.08$	$92.5 \pm 0.12$	$96.3 \pm 0.20$	$69.2 \pm 0.67$	$80.3 \pm 0.68$	$58.1 \pm 0.06$	$73.1 \pm 0.06$
AdaptFormer	$87.9 \pm 0.10$	$94.8 \pm 0.21$	$86.2 \pm 0.19$	$87.6 \pm 0.24$	$93.2 \pm 0.15$	$93.3 \pm 0.68$	$97.0 \pm 0.81$	$75.0 \pm 0.11$	$84.8 \pm 0.08$	$61.1 \pm 0.15$	$75.5 \pm 0.12$
LoRA	$87.7 \pm 0.59$	$94.6 \pm 0.50$	$85.1 \pm 0.64$	$87.8 \pm 0.24$	$93.3 \pm 0.13$	$93.0 \pm 0.14$	$96.6 \pm 0.11$	$71.4 \pm 0.54$	$82.1 \pm 0.62$	$59.0 \pm 0.19$	$74.0 \pm 0.17$
Adapter	$88.2 \pm 0.44$	$94.8 \pm 0.34$	$86.7 \pm 0.92$	$87.7 \pm 0.23$	$93.2 \pm 0.16$	$93.4 \pm 0.12$	$97.1 \pm 0.15$	$74.4 \pm 0.16$	$84.3 \pm 0.28$	$60.5 \pm 0.10$	$75.1 \pm 0.08$
HQ-SAM	$85.1 \pm 0.10$	$92.6 \pm 0.10$	$81.0 \pm 0.61$	$86.3 \pm 0.32$	$92.4 \pm 0.19$	$91.1 \pm 0.50$	$95.5 \pm 0.57$	$66.2 \pm 0.44$	$77.8 \pm 0.43$	$54.9 \pm 0.16$	$70.6 \pm 0.13$
SU-SAM	$88.3 \pm 0.21$	$95.0 \pm 0.22$	$86.2 \pm 0.59$	$87.8 \pm 0.18$	$93.2 \pm 0.09$	$93.8 \pm 0.02$	$97.5 \pm 0.06$	$74.7 \pm 0.53$	$84.5 \pm 0.56$	$60.2 \pm 0.26$	$74.8 \pm 0.22$
ConvLoRA-SAM	$87.5 \pm 0.39$	$94.5 \pm 0.17$	$85.4 \pm 0.41$	$87.7 \pm 0.22$	$93.2 \pm 0.11$	$92.9 \pm 0.13$	$96.6 \pm 0.28$	$71.4 \pm 0.44$	$82.2 \pm 0.37$	$59.6 \pm 0.22$	$74.4 \pm 0.20$
LoRA+Ours	$88.3 \pm 0.05$	$95.2 \pm 0.00$	$85.8 \pm 0.59$	$88.1 \pm 0.08$	$93.5 \pm 0.05$	$93.4 \pm 0.11$	$96.8 \pm 0.09$	$72.2 \pm 0.06$	$82.8 \pm 0.04$	$59.9 \pm 0.20$	$74.6 \pm 0.17$
Adapter+Ours	$88.6 \pm 0.09$	$95.1 \pm 0.05$	$87.1 \pm 0.37$	$88.0 \pm 0.05$	$93.4 \pm 0.00$	$94.4 \pm 0.12$	$97.9 \pm 0.09$	$75.6 \pm 0.27$	$85.2 \pm 0.23$	$61.4 \pm 0.30$	$75.8 \pm 0.27$

InfoSAM outperforms other PEFT techniques across various datasets from different domains.



# Experiments



## □ Compare with distillation baselines across various domains

METHOD	NATURAL IMAGES			MEDICAL				AGRICULTURE		REMOTE SENSING	
	CAMO			ISIC 2017		Kvasir		Leaf		Road	
	$S_\alpha \uparrow$	$E_\phi \uparrow$	$F_\beta^\omega \uparrow$	Jac $\uparrow$	Dice $\uparrow$	$S_\alpha \uparrow$	$E_\phi \uparrow$	IoU $\uparrow$	Dice $\uparrow$	IoU $\uparrow$	Dice $\uparrow$
Teacher	$79.7 \pm 0.02$	$88.8 \pm 0.09$	$79.6 \pm 0.01$	$61.0 \pm 0.12$	$71.7 \pm 0.14$	$83.0 \pm 0.10$	$88.8 \pm 0.29$	$37.6 \pm 0.11$	$47.0 \pm 0.16$	$7.2 \pm 0.24$	$12.9 \pm 0.29$
Student	$88.2 \pm 0.44$	$94.8 \pm 0.34$	$86.7 \pm 0.92$	$87.7 \pm 0.23$	$93.2 \pm 0.16$	$93.4 \pm 0.12$	$97.1 \pm 0.15$	$74.4 \pm 0.16$	$84.3 \pm 0.28$	$60.5 \pm 0.10$	$75.1 \pm 0.08$
Logit	$88.4 \pm 0.08$	$94.9 \pm 0.05$	$87.1 \pm 0.22$	$87.2 \pm 0.43$	$92.9 \pm 0.29$	$93.2 \pm 0.19$	$96.5 \pm 0.19$	$73.0 \pm 0.35$	$83.3 \pm 0.29$	$50.9 \pm 0.08$	$67.2 \pm 0.06$
PKD	$87.0 \pm 0.43$	$94.1 \pm 0.23$	$84.3 \pm 0.97$	$86.5 \pm 0.26$	$92.5 \pm 0.17$	$92.2 \pm 0.25$	$96.0 \pm 0.17$	$70.2 \pm 1.15$	$81.1 \pm 1.08$	$56.9 \pm 0.61$	$72.2 \pm 0.56$
PKT	$87.8 \pm 0.40$	$94.5 \pm 0.35$	$86.2 \pm 0.46$	$87.4 \pm 0.12$	$93.0 \pm 0.07$	$93.7 \pm 0.41$	$97.3 \pm 0.53$	$74.2 \pm 0.51$	$84.2 \pm 0.52$	$60.7 \pm 0.20$	$75.2 \pm 0.16$
IBD	$85.2 \pm 0.47$	$92.6 \pm 0.35$	$82.4 \pm 0.31$	$85.1 \pm 0.74$	$91.7 \pm 0.45$	$91.5 \pm 0.14$	$95.3 \pm 0.05$	$72.2 \pm 0.12$	$82.7 \pm 0.07$	$44.9 \pm 0.18$	$61.5 \pm 0.18$
VID	$87.9 \pm 0.22$	$94.8 \pm 0.34$	$86.3 \pm 0.32$	$87.6 \pm 0.44$	$93.1 \pm 0.29$	$93.7 \pm 0.16$	$97.4 \pm 0.07$	$75.1 \pm 0.08$	$84.9 \pm 0.17$	$60.7 \pm 0.19$	$75.4 \pm 0.19$
SemCKD	$86.2 \pm 0.16$	$93.5 \pm 0.21$	$82.8 \pm 1.54$	$85.4 \pm 0.27$	$91.8 \pm 0.19$	$92.4 \pm 0.07$	$96.2 \pm 0.03$	$72.0 \pm 0.04$	$82.8 \pm 0.10$	$53.5 \pm 0.17$	$69.4 \pm 0.17$
ReviewKD	$86.7 \pm 0.07$	$94.0 \pm 0.09$	$84.6 \pm 0.63$	$85.5 \pm 0.26$	$91.9 \pm 0.15$	$92.4 \pm 0.33$	$96.4 \pm 0.26$	$72.6 \pm 0.64$	$83.1 \pm 0.47$	$57.3 \pm 0.11$	$72.6 \pm 0.11$
TinySAM	$83.7 \pm 0.39$	$91.6 \pm 0.31$	$81.1 \pm 0.35$	$79.4 \pm 1.12$	$87.8 \pm 0.84$	$88.5 \pm 0.31$	$93.5 \pm 0.24$	$48.6 \pm 1.14$	$61.0 \pm 0.95$	$25.7 \pm 1.19$	$39.6 \pm 1.71$
MobileSAM	$87.1 \pm 0.36$	$94.1 \pm 0.27$	$85.1 \pm 0.09$	$86.7 \pm 0.13$	$92.6 \pm 0.09$	$92.5 \pm 0.12$	$96.3 \pm 0.14$	$71.9 \pm 0.30$	$82.6 \pm 0.39$	$59.2 \pm 0.09$	$74.1 \pm 0.08$
InfoSAM(Ours)	<b><math>88.6 \pm 0.09</math></b>	<b><math>95.1 \pm 0.05</math></b>	<b><math>87.1 \pm 0.37</math></b>	<b><math>88.0 \pm 0.05</math></b>	<b><math>93.4 \pm 0.00</math></b>	<b><math>94.4 \pm 0.12</math></b>	<b><math>97.9 \pm 0.09</math></b>	<b><math>75.6 \pm 0.27</math></b>	<b><math>85.2 \pm 0.23</math></b>	<b><math>61.4 \pm 0.30</math></b>	<b><math>75.8 \pm 0.27</math></b>

Most distillation methods harm PEFT due to the weak teacher, often underperforming vanilla fine-tuning. In contrast, InfoSAM distills only essential knowledge from the teacher.

# Experiments



## □ Extended Experiment with SAM2

(a) PEFT Methods Comparison

METHOD	MEDICAL	AGRICULTURE	REMOTE SENSING
	$S_{\alpha}$ (Kvasir)	IoU (Leaf)	IoU (Road)
SAM2	$87.1 \pm 0.12$	$42.7 \pm 0.32$	$6.9 \pm 0.13$
decoder-only	$93.2 \pm 0.07$	$71.8 \pm 0.58$	$48.5 \pm 0.47$
BitFit	$93.8 \pm 0.09$	$75.4 \pm 0.29$	$59.2 \pm 0.26$
AdaptFormer	$93.7 \pm 0.19$	$73.6 \pm 1.10$	$59.9 \pm 0.35$
LoRA	$93.7 \pm 0.10$	$75.9 \pm 0.40$	$60.8 \pm 0.32$
Adapter	$94.4 \pm 0.06$	$76.8 \pm 0.56$	$60.9 \pm 0.14$
<b>LoRA+Ours</b>	<b><math>94.0 \pm 0.09</math></b>	<b><math>76.1 \pm 0.38</math></b>	<b><math>60.9 \pm 0.05</math></b>
<b>Adapter+Ours</b>	<b><math>94.5 \pm 0.17</math></b>	<b><math>77.3 \pm 0.14</math></b>	<b><math>61.3 \pm 0.05</math></b>

(b) Distillation Methods Comparison

METHOD	MEDICAL	AGRICULTURE	REMOTE SENSING
	$S_{\alpha}$ (Kvasir)	IoU (Leaf)	IoU (Road)
Teacher	$87.1 \pm 0.12$	$42.7 \pm 0.32$	$6.9 \pm 0.13$
Student	$94.4 \pm 0.06$	$76.8 \pm 0.56$	$60.9 \pm 0.14$
PKT	$94.0 \pm 0.25$	$74.8 \pm 0.14$	$57.3 \pm 0.07$
VID	$94.1 \pm 0.47$	$77.2 \pm 0.37$	$61.1 \pm 0.38$
ReviewKD	$93.4 \pm 0.10$	$72.7 \pm 0.37$	$55.9 \pm 0.50$
TinySAM	$89.4 \pm 0.10$	$45.2 \pm 0.76$	$23.9 \pm 2.61$
MobileSAM	$93.3 \pm 0.15$	$74.1 \pm 0.35$	$52.3 \pm 0.46$
<b>InfoSAM2(Ours)</b>	<b><math>94.5 \pm 0.17</math></b>	<b><math>77.3 \pm 0.14</math></b>	<b><math>61.3 \pm 0.05</math></b>

InfoSAM consistently performs well with SAM2, thanks to its structure-independent, information-theoretic foundation.

# Experiments



## ❑ Ablation Study

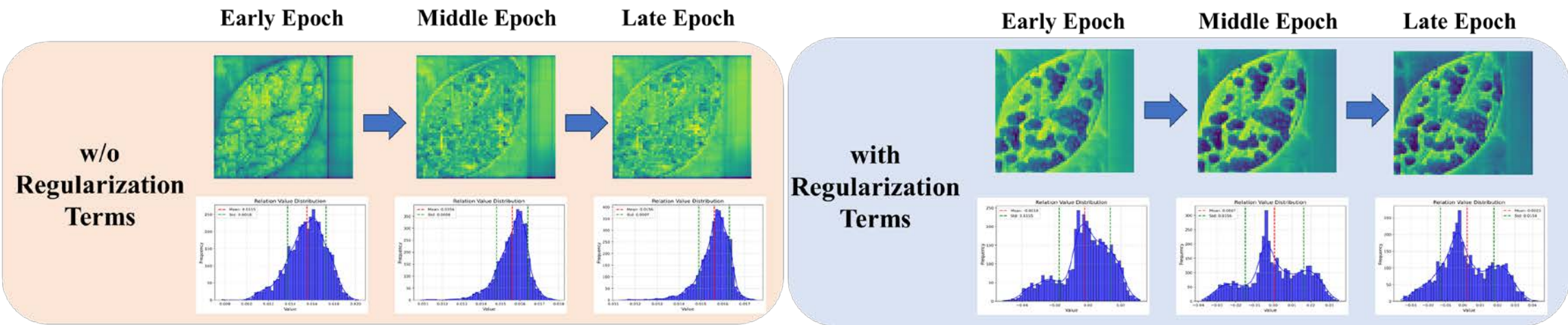
✓ Ablation study results of two losses

$L_r$	$L_d$	MEDICAL	AGRICULTURE	REMOTE SENSING
		$S_\alpha$ (Kvasir)	IoU (Leaf)	IoU (Road)
		93.4	74.4	60.5
	✓	93.6 (+0.2)	75.2 (+0.8)	61.0 (+0.5)
✓	✓	94.4 (+1.0)	75.6 (+1.2)	61.4 (+0.9)

✓ Effects of the Relation Module

MODEL	METHOD	AGRICULTURE	REMOTE SENSING
		IoU (Leaf)	IoU (Road)
TinySAM	w/o RM	$48.6 \pm 1.14$	$28.7 \pm 1.69$
	w RM	<b><math>50.3 \pm 0.76</math></b>	<b><math>33.9 \pm 0.32</math></b>
MobileSAM	w/o RM	$71.9 \pm 0.30$	$59.2 \pm 0.09$
	w RM	<b><math>73.8 \pm 0.22</math></b>	<b><math>61.3 \pm 0.35</math></b>

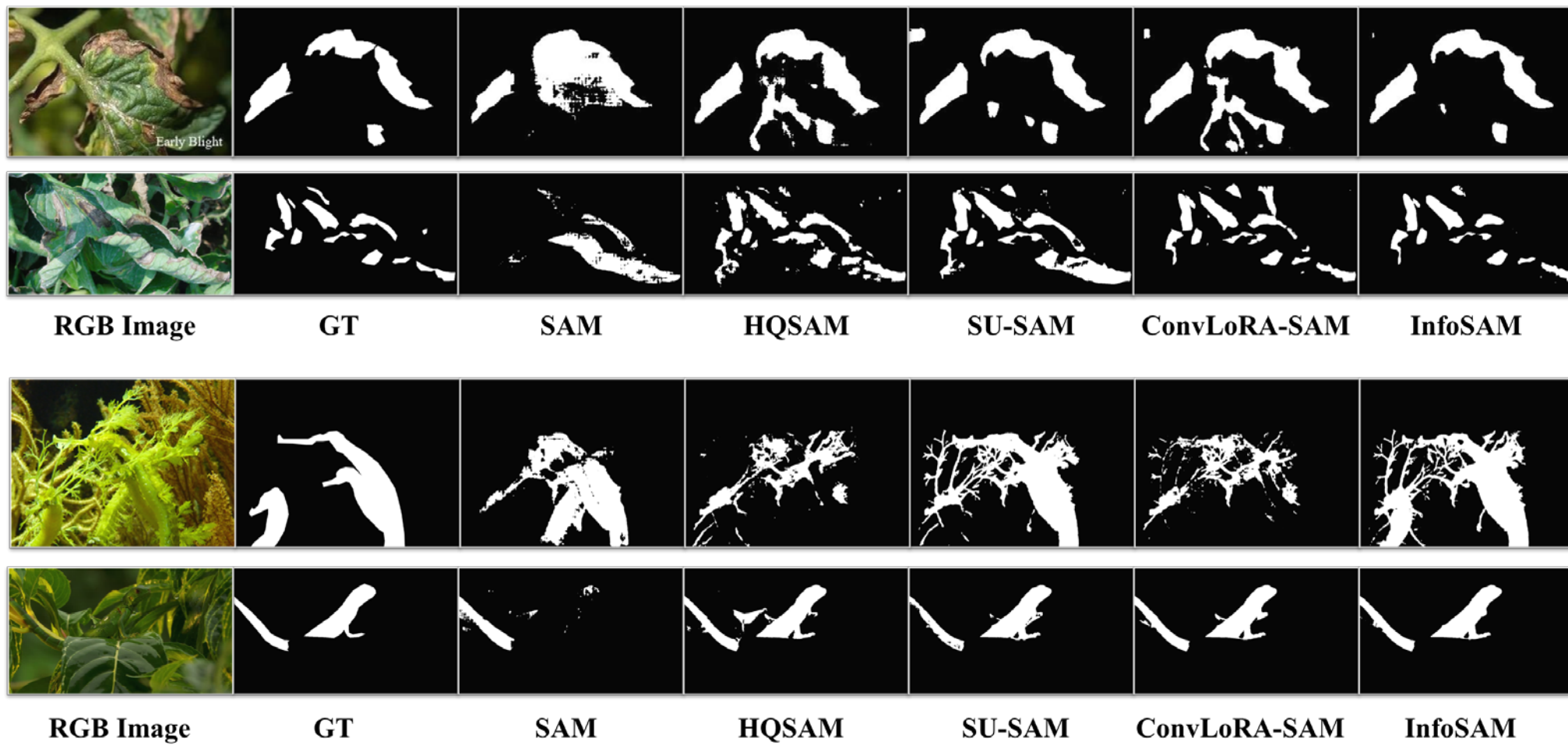
✓ Evolution of relation maps and their statistical distributions over epochs, without and with the regularization term.



# Experiments



## □ Visualization results



# Reference



- [1] Kirillov A, Mintun E, Ravi N, et al. Segment anything[C]//Proceedings of the IEEE/CVF international conference on computer vision. 2023: 4015-4026
- [2] Peng Z, Xu Z, Zeng Z, et al. Parameter efficient fine-tuning via cross block orchestration for segment anything model[C]//Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2024: 3743-3752.
- [3] Ke L, Ye M, Danelljan M, et al. Segment anything in high quality[J]. Advances in Neural Information Processing Systems, 2023, 36: 29914-29934.
- [4] Xiao A, Xuan W, Qi H, et al. CAT-SAM: conditional tuning network for few-shot adaptation of segmentation anything model[J]. arXiv e-prints, 2024: arXiv: 2402.03631.