



# CONVOLUTION MEETS LoRA: PARAMETER EFFICIENT FINETUNING FOR SEGMENT ANYTHING MODEL *[ICLR 24]*

Zihan Zhong, Zhiqiang Tang, Tong He, Haoyang Fang, Chun Yuan  
Tsinghua University, Amazon Web Services

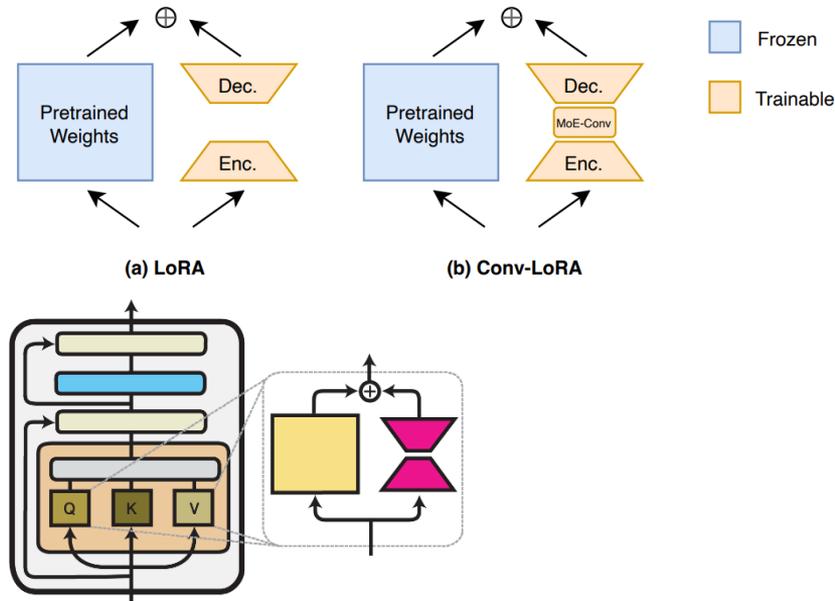
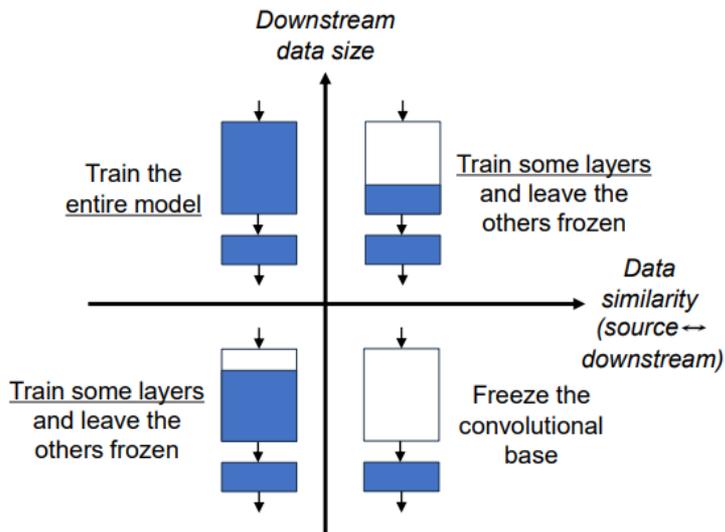
Sunghyun Ahn  
[skd@yonsei.ac.kr](mailto:skd@yonsei.ac.kr)

<2024/06/05>

# 1 Introduction

## Parameter Efficient FineTuning

- ➡ PEFT: pre-trained 모델의 모든 parameter를 학습시키지 않고, 적은 양의 parameter만 학습하여 사전 지식을 유지하는 기법
- ➡ LoRA: pre-trained 모델의 weight를 update하지 않고도 fully-tuning한 결과와 비슷하거나 더 좋은 성능을 보인 PEFT 방법
- ➡ ConvLoRA: LoRA에 Convolution block을 추가하여 vision-specific inductive bias를 추가한 방법



# 1 Introduction

## LoRA [ICLR 21]

- Low Rank Decomposition: 큰 차원의 Matrix 연산을 낮은 차원으로 분해하는 것
- Training: 사전 학습한 모델의 weight들은 업데이트 하지 않고, LoRA의 rank decomposition matrices의 weight들만 업데이트 (기존 크기의 0.01%)
- Testing: pretrained weights와 trainable weights를 더하므로, 추가적인 Inference Latency가 발생하지 않음 ( $W = W_0 + BA$ )

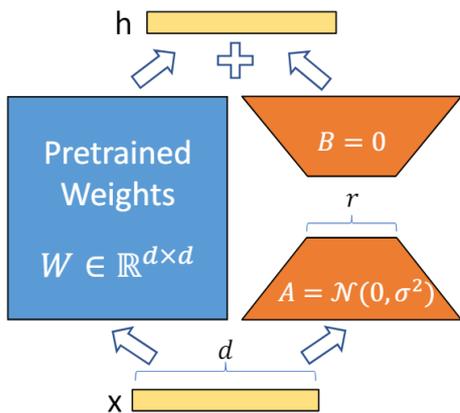


Figure 1: Our reparametrization. We only train  $A$  and  $B$ .

$$h = W_0x + \Delta Wx = W_0x + BAx$$

$$W_0 \in \mathbb{R}^{d \times k}, B \in \mathbb{R}^{d \times r}, A \in \mathbb{R}^{r \times k}, r \ll \min(d, k)$$

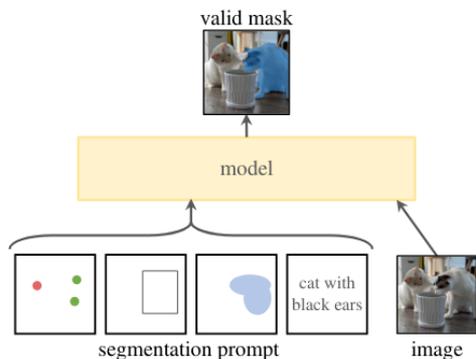
$$\max_{\Phi} \sum_{(x,y) \in \mathcal{Z}} \sum_{t=1}^{|y|} \log(P_{\Phi}(y_t|x, y_{<t}))$$

$$\max_{\Theta} \sum_{(x,y) \in \mathcal{Z}} \sum_{t=1}^{|y|} \log(p_{\Phi_0 + \Delta\Phi(\Theta)}(y_t|x, y_{<t}))$$

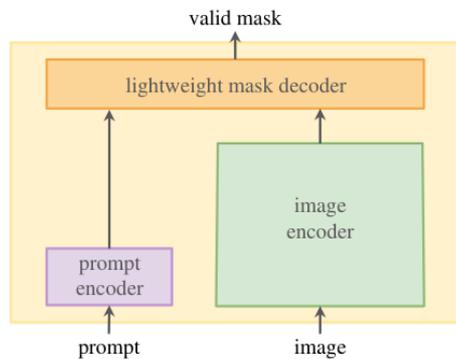
# 1 Introduction

## SAM [ICCV 23]

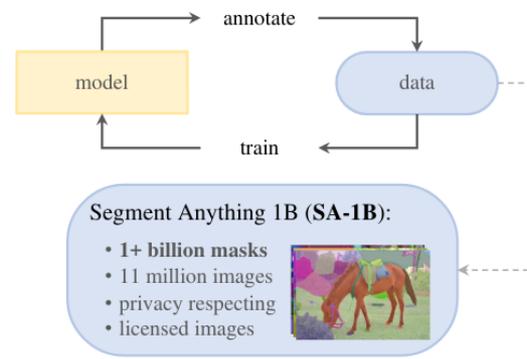
- ➡ **Task:** zero-shot segmentation을 위한 promptable segmentation 제안 (image + points/bbox/mask/text -> valid mask)
- ➡ **Model:** prompt와 image를 입력으로 받아 실시간으로 mask를 예측하는 모델 제안 (encoder: pretrained MAE, CLIP)
- ➡ **Data:** SA-1B (Large Instance Segmentation Dataset) 제공 (human-annotated data + model-annotated data)



(a) **Task:** promptable segmentation



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



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

Figure 1: We aim to build a foundation model for segmentation by introducing three interconnected components: a promptable segmentation *task*, a segmentation *model* (SAM) that powers data annotation and enables zero-shot transfer to a range of tasks via prompt engineering, and a *data* engine for collecting SA-1B, our dataset of over 1 billion masks.

# 1 Introduction

## 🔗 Conv-LoRA

🔗 **Key 1:** 기존 LoRA layer 사이에 lightweight한 convolution block 삽입

Convolution의 고유 특징인 vision-specific inductive bias를 활용 가능

🔗 **Key 2:** MoE(Mixture-of-Experts) 컨셉 구현

여러 experts를 두고, 각 expert가 다른 scale의 image feature를 확인 ->  
다양한 scale의 이미지 특징을 추출하므로, 다양한 객체를 탐지 가능

🔗 **Key 3:** SAM의 decoder에 다중 클래스 예측 MLP 추가

SAM을 통해 Semantic Segmentation을 가능하게 함

## 2 Method

### Conv-LoRA

- LoRA의 중간부에 Convolution operations 삽입
- 단순한 Conv block 대신, MoE-Conv block을 제안

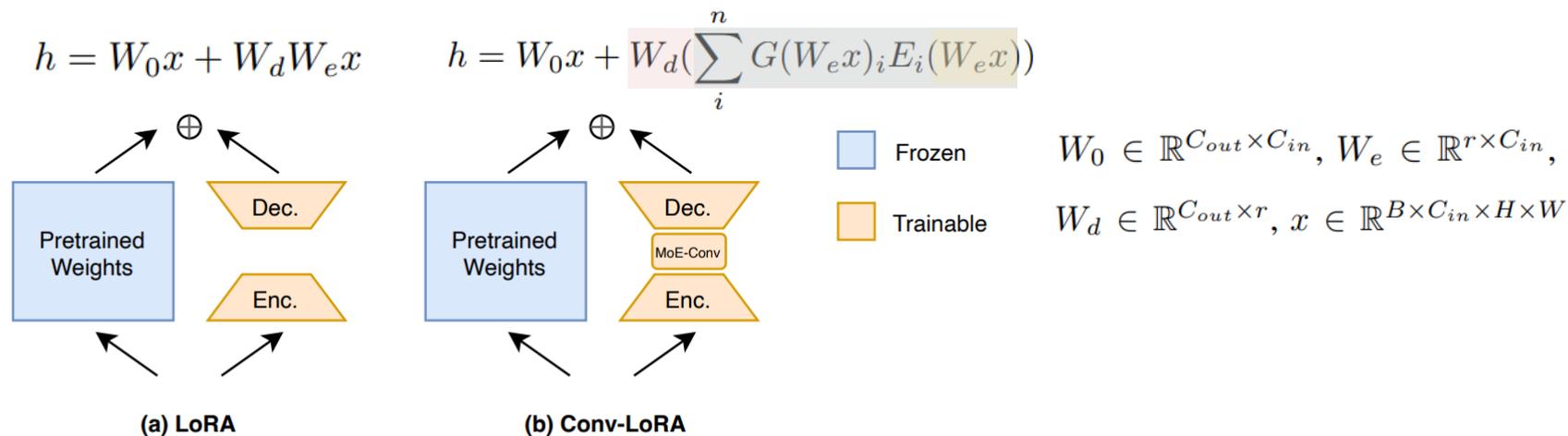
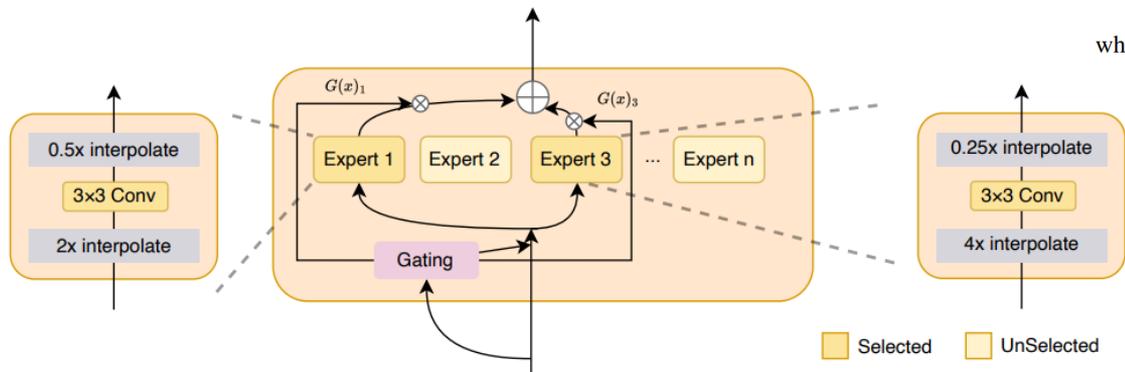


Figure 2: LoRA vs. Conv-LoRA. Both LoRA and Conv-LoRA add an extra trainable encoder-decoder structure parallel to the frozen pre-trained weights. Inside the bottleneck of LoRA, Conv-LoRA inserts lightweight convolution operations managed by MoE with negligible extra parameters.

## 2 Method

### MoE-Conv

- ↔ N개의 expert networks와 gating module로 구성됨 (expert: 다른 스케일의 feature 확인, gating: dynamic하게 expert 선택)
- ↔ 다양한 scale의 이미지 특징을 추출하므로, 다양한 크기의 객체를 탐지할 수 있음
- ↔ Multi-scale 방식(i.e. Swin)에 비해 computational cost가 적어 효율적임



$$E_i(x) = \text{Interpolate}(\text{Conv}_{3 \times 3}(\text{Interpolate}(x, s_i)), 1/s_i) \quad (3)$$

For instance, if  $s_i = 4$ , expert  $E_i$  would initially upscale the feature maps by a factor of 4x, apply the  $\text{Conv}_{3 \times 3}$  operation, and finally, downscale the feature maps by 4x.

$$x_h = \text{Reshape}(\text{AvgPool}(x), (B, r))$$

$$H(x)_i = (x_h \cdot W_g)_i + \text{StandardNormal}() \cdot \text{Softplus}((x_h \cdot W_{noise})_i)$$

$$G(x) = \text{Softmax}(\text{KeepTopK}(H(x), k))$$

$$\text{where } \text{KeepTopK}(v, k)_i = \begin{cases} v_i & \text{if } v_i \text{ is in the top } k \text{ elements of } v. \\ -\infty & \text{otherwise.} \end{cases}$$

$$h = W_0 x + W_d \left( \sum_i^n G(W_e x)_i E_i(W_e x) \right)$$

$$x_h \in \mathbb{R}^{B \times r} \quad W_g \in \mathbb{R}^{r \times n}$$

## 2 Method

### Multi-class Segmentation with SAM

- Mask-based method (i.e. Maskformer)를 따라  $N$ 개의 binary masks와  $N$ 개의 classification scores를 예측함
- Multi-class 예측을 위해 prompt encoder를 제거 후 lightweight한 MLP를 추가함
- SAM의 구조를 변경하여 end-to-end로 학습이 가능해짐

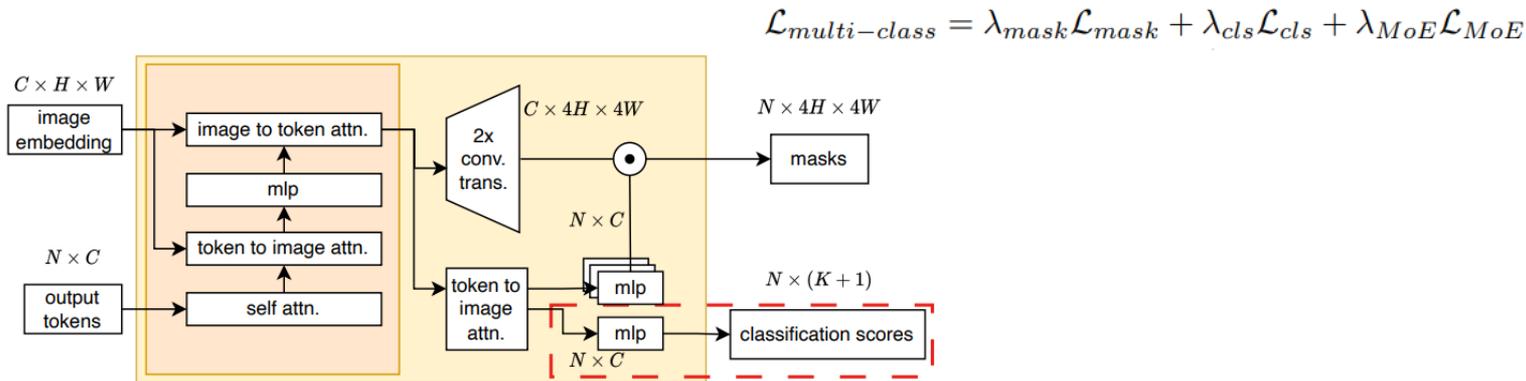


Figure 4: The modified SAM's mask decoder for multi-class semantic segmentation. The classification module (within the red dashed box) is new added compared to original SAM's mask decoder.  $N$  is the number of output mask tokens,  $K$  is the number of classes,  $C$  is the number of channels,  $H$  and  $W$  indicate the height and width of the feature map (we omit 'batch size' for simplicity).

# 3 Experiments

## Binary class semantic segmentation

- Medical, Natural, Agriculture, Remote sensing 환경에서 실험 진행
- SAM trained from scratch 성능이 가장 안 좋은 이유는 PEFT 방법들과 달리 SAM의 사전 지식을 잃어버렸기 때문
- Vision-specific한 inductive bias가 추가됨으로써 가장 우수한 성능을 보임

Method	#Params (M) / Ratio (%)	Medical						Natural Images			Agriculture		Remote Sensing		
		Kvasir		CVC-612		ISIC 2017		CAMO			SBU	Leaf		Road	
		$S_\alpha \uparrow$	$E_\phi \uparrow$	$S_\alpha \uparrow$	$E_\phi \uparrow$	Jac $\uparrow$	Dice $\uparrow$	$S_\alpha \uparrow$	$E_\phi \uparrow$	$F_\beta^\omega \uparrow$	BER $\downarrow$	IoU $\uparrow$	Dice $\uparrow$	IoU $\uparrow$	Dice $\uparrow$
<i>Domain Specific</i>	* / 100%	90.9	94.4	<u>92.6</u>	<u>95.5</u>	<u>80.1</u>	<u>87.5</u>	80.8	85.8	73.1	3.56	62.3	74.1	59.1	73.0
<i>SAM trained from scratch</i>	641.09 / 100%	78.5	82.4	85.9	91.6	73.8	82.5	61.9	67.0	40.5	5.53	52.1	65.5	55.6	71.1
decoder-only	3.51 / 0.55%	86.5	89.5	85.5	89.9	69.7	79.5	78.5	83.1	69.8	14.58	50.8	63.8	48.6	65.1
BitFit	3.96 / 0.62%	90.8 $\pm$ 0.57	93.8 $\pm$ 0.98	89.0 $\pm$ 0.40	91.6 $\pm$ 0.90	76.4 $\pm$ 0.45	84.7 $\pm$ 0.35	86.8 $\pm$ 0.33	90.7 $\pm$ 0.28	81.5 $\pm$ 0.19	3.16 $\pm$ 0.128	71.4 $\pm$ 1.15	81.7 $\pm$ 1.01	60.6 $\pm$ 0.15	75.2 $\pm$ 0.11
Adapter	3.92 / 0.61%	91.2 $\pm$ 0.23	94.0 $\pm$ 0.16	89.3 $\pm$ 0.43	92.0 $\pm$ 0.63	76.7 $\pm$ 0.66	85.0 $\pm$ 0.56	87.7 $\pm$ 0.10	91.3 $\pm$ 0.40	82.8 $\pm$ 0.35	2.84 $\pm$ 0.093	72.1 $\pm$ 0.47	82.4 $\pm$ 0.36	61.5 $\pm$ 0.11	75.9 $\pm$ 0.12
VPT	4.00 / 0.62%	91.5 $\pm$ 0.23	94.3 $\pm$ 0.06	91.0 $\pm$ 0.94	93.7 $\pm$ 1.41	76.9 $\pm$ 0.94	85.1 $\pm$ 0.75	87.4 $\pm$ 0.66	91.4 $\pm$ 0.68	82.1 $\pm$ 0.75	2.70 $\pm$ 0.055	73.6 $\pm$ 0.28	83.8 $\pm$ 0.26	60.2 $\pm$ 1.87	74.9 $\pm$ 1.50
LST	11.49 / 1.77%	89.7 $\pm$ 0.25	93.3 $\pm$ 0.37	89.4 $\pm$ 0.37	92.4 $\pm$ 0.54	76.4 $\pm$ 1.05	84.9 $\pm$ 0.79	83.3 $\pm$ 0.28	88.0 $\pm$ 0.23	77.1 $\pm$ 0.012	3.18 $\pm$ 0.012	70.2 $\pm$ 0.87	81.1 $\pm$ 0.82	60.2 $\pm$ 0.26	74.9 $\pm$ 0.22
SAM-Adapter	3.98 / 0.62%	89.6 $\pm$ 0.24	92.5 $\pm$ 0.10	89.6 $\pm$ 0.22	92.4 $\pm$ 1.09	76.1 $\pm$ 0.45	84.6 $\pm$ 0.37	85.6 $\pm$ 0.28	89.6 $\pm$ 0.55	79.8 $\pm$ 0.89	3.14 $\pm$ 0.063	71.4 $\pm$ 0.20	82.1 $\pm$ 0.10	60.6 $\pm$ 0.06	75.2 $\pm$ 0.04
SSF	4.42 / 0.69%	91.3 $\pm$ 0.87	93.9 $\pm$ 1.49	89.6 $\pm$ 0.37	91.9 $\pm$ 0.79	76.6 $\pm$ 0.19	85.0 $\pm$ 0.14	87.5 $\pm$ 0.11	91.4 $\pm$ 0.16	82.6 $\pm$ 0.12	3.19 $\pm$ 0.044	71.5 $\pm$ 0.63	81.8 $\pm$ 0.44	61.6 $\pm$ 0.03	76.0 $\pm$ 0.02
LoRA	4.00 / 0.62%	91.2 $\pm$ 0.28	93.8 $\pm$ 0.22	90.7 $\pm$ 0.04	92.5 $\pm$ 0.41	76.6 $\pm$ 0.23	84.9 $\pm$ 0.22	88.0 $\pm$ 0.24	91.9 $\pm$ 0.42	82.8 $\pm$ 0.16	2.74 $\pm$ 0.079	73.7 $\pm$ 0.20	83.6 $\pm$ 0.13	62.2 $\pm$ 0.21	76.5 $\pm$ 0.18
Conv-LoRA	4.02 / 0.63%	<b>92.0</b> $\pm$ 0.15	<b>94.7</b> $\pm$ 0.16	<b>91.3</b> $\pm$ 0.69	<b>94.0</b> $\pm$ 0.79	<b>77.6</b> $\pm$ 0.57	<b>85.7</b> $\pm$ 0.36	<b>88.3</b> $\pm$ 0.46	<b>92.4</b> $\pm$ 0.31	<b>84.0</b> $\pm$ 0.34	<b>2.54</b> $\pm$ 0.081	<b>74.5</b> $\pm$ 0.39	<b>84.3</b> $\pm$ 0.34	<b>62.6</b> $\pm$ 0.36	<b>76.8</b> $\pm$ 0.27

Table 1: Results on binary semantic segmentation. ‘# Params (M) / Ratio (%)’ represents the number of trainable parameters and its proportion relative to the total number. *Domain Specific* is a placeholder referring to methods that specifically designed for the tasks. ‘Underlined’ denotes the better results compared to PEFT methods. See appendix D for more details. Compared to LoRA, Conv-LoRA incurs negligible parameter overhead, but delivers a clear performance boost.

# 3 Experiments

## Multi class semantic segmentation

- 다양한 클래스를 지닌 multi-class dataset에서도 가장 좋은 성능을 보임
- Decoder만 학습하면 Encoder가 high-level 정보를 추출하지 못하지만, Encoder를 PEFT로 학습하면 성능이 크게 향상됨

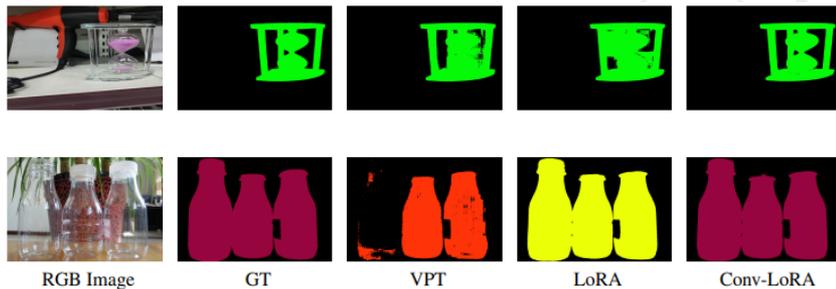


Figure 14: Visualization Results on Multi-class Transparent Object Segmentation.

Method	#Params (M) / Ratio (%)	Easy				Hard			
		Acc ↑	mIoU ↑	MAE ↓	MBER ↓	Acc ↑	mIoU ↑	MAE ↓	MBER ↓
<i>TransLab</i>	42.19 / 100%	95.77	92.23	0.036	3.12	83.04	72.10	0.166	13.30
decoder-only	3.51 / 0.55%	94.68	88.54	0.050	4.24	83.53	68.30	0.186	14.37
VPT	4.00 / 0.62%	98.31	95.73	0.017	1.52	90.42	83.38	0.083	7.21
LoRA	4.00 / 0.62%	98.44	96.26	0.016	1.35	91.94	83.95	0.083	6.35
Conv-LoRA	4.02 / 0.63%	<b>98.63</b>	<b>96.45</b>	<b>0.015</b>	<b>1.27</b>	<b>93.05</b>	<b>84.37</b>	<b>0.075</b>	<b>6.25</b>

Method	# Params (M) / Ratio (%)	Acc ↑	mIoU ↑	Category IoU ↑											
				bg	shelf	jar	freezer	window	door	eyeglass	cup	wall	bowl	bottle	box
<i>TransLab</i>	42.19 / 100%	92.67	69.00	93.90	54.36	64.48	65.14	54.58	57.72	79.85	81.61	72.82	69.63	77.50	56.43
<i>Trans2Seg</i>	56.20 / 100%	94.14	<u>72.15</u>	95.35	<u>53.43</u>	<u>67.82</u>	64.20	<u>59.64</u>	60.56	<u>88.52</u>	<u>86.67</u>	75.99	<u>73.98</u>	<u>82.43</u>	<u>57.17</u>
decoder-only	3.51 / 0.55%	90.66	49.97	93.66	32.75	39.96	35.87	50.70	45.89	57.38	73.16	69.36	54.23	56.58	33.77
VPT	4.00 / 0.62%	94.42	62.81	97.41	29.76	52.82	62.09	55.54	63.61	81.12	83.40	79.61	65.29	72.92	44.77
LoRA	4.00 / 0.62%	94.80	66.01	97.50	42.17	57.82	<b>64.35</b>	53.44	64.08	<b>87.28</b>	<b>85.28</b>	80.43	63.67	77.97	49.56
Conv-LoRA	4.02 / 0.63%	<b>95.07</b>	<b>67.09</b>	<b>97.66</b>	<b>50.51</b>	<b>58.44</b>	51.70	<b>55.69</b>	<b>65.22</b>	85.23	84.84	<b>80.97</b>	<b>72.84</b>	<b>79.83</b>	<b>52.73</b>

Table 2: Results on multi-class semantic segmentation. The tables are the results for three-class and twelve-class semantic segmentation respectively.

# 3 Experiments

## Mean attention distance

- MAE로 학습된 ViT는 후반부에서 Mean Attention Distance가 길게 나타남 (후반부로 갈수록 fine-grained information이 부족함)
- ConvLoRA로 학습된 SAM의 Encoder는 후반부에서 Mean Attention Distance가 짧고 나타남 (fine-grained detail이 풍부하여 segmentation에 탁월함)

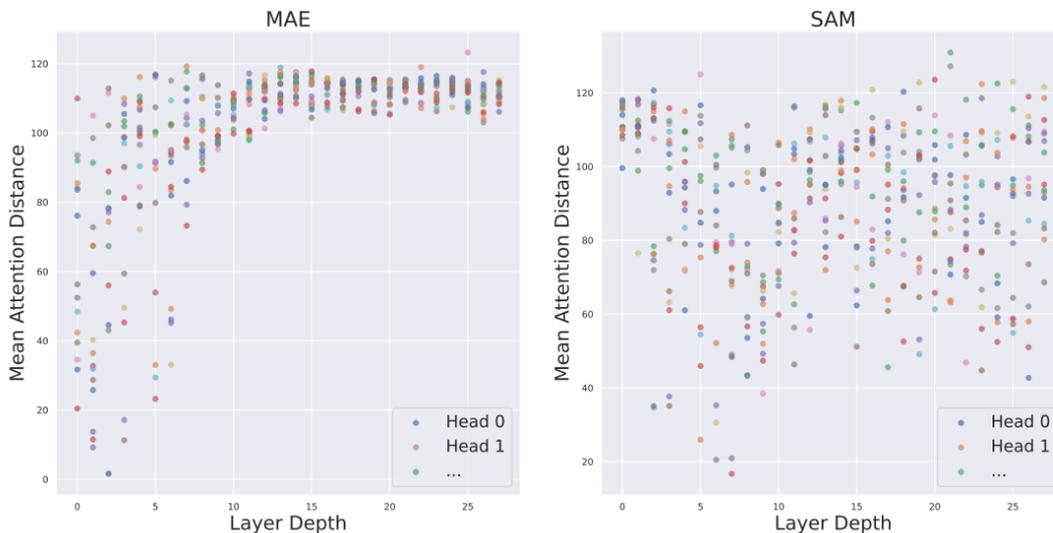


Figure 5: Mean attention distance of each attention head, with each dot indicating the mean distance across images for one of the 16 heads at one layer. In contrast to MAE, SAM retains the ability to incorporate local information even in deeper layers.

# 3 Experiments

## Further Analysis

Method	ISIC 2017				Method	Training Speed (Iter / s) $\uparrow$	Training Memory (GB) $\downarrow$
	Jac $\uparrow$	T-Jac $\uparrow$	Dice $\uparrow$	Acc $\uparrow$			
Multi-scale	77.4	69.5	85.4	93.7		0.79	23.4
MoE	<b>77.9</b>	<b>70.3</b>	<b>85.9</b>	<b>93.9</b>		<b>1.22</b>	<b>21.7</b>

Table 3: MoE vs. Multi-scale, the latter fuses features from all scales simultaneously. The performance and the training cost illustrate the effectiveness and efficiency of dynamic selecting the experts, i.e., injecting the local prior into the feature maps of different scales dynamically.

LoRA	#Params (M)	Test	
		Acc $\uparrow$	mIoU $\uparrow$
$r = 3$	4.00	94.80	66.01
$r = 6$	4.49	95.15	66.24
$r = 12$	5.48	95.23	66.69
$r = 24$	7.44	<b>95.14</b>	<b>67.02</b>

Table 11: The performance trend when increasing the rank  $r$  of LoRA for twelve-class segmentation.

Method	# Params / Ratio(%)	CVC-ColonDB				ETIS			
		$S_\alpha \uparrow$	$E_\phi \uparrow$	$F_\beta^\omega \uparrow$	MAE $\downarrow$	$S_\alpha \uparrow$	$E_\phi \uparrow$	$F_\beta^\omega \uparrow$	MAE $\downarrow$
VPT	4.00 / 0.62%	85.0	88.4	75.0	3.6	86.5	87.4	70.1	1.9
Conv-LoRA	4.03 / 0.63%	85.1	88.8	75.7	<b>3.4</b>	86.8	<b>88.5</b>	70.2	<b>1.7</b>
Conv-LoRA+VPT	4.23 / 0.66 %	<b>86.0</b>	<b>89.0</b>	<b>76.8</b>	3.5	<b>87.7</b>	88.1	<b>70.3</b>	<b>1.7</b>

Table 12: Performance of combining Conv-LoRA and VPT on Polyp Segmentation.

Method	# Params / Ratio(%)	Leaf			ISIC 2017			
		IoU $\uparrow$	Dice $\uparrow$	Acc $\uparrow$	Jac $\uparrow$	T-Jac $\uparrow$	Dice $\uparrow$	Acc $\uparrow$
Inception	4.01 / 0.63%	72.1	82.3	95.0	76.9	68.5	84.9	93.3
MoE-Conv	4.01 / 0.63%	<b>74.0</b>	<b>83.8</b>	<b>95.6</b>	<b>77.4</b>	<b>69.5</b>	<b>85.5</b>	<b>93.9</b>

Table 13: MoE-Conv vs. Blocks with Various Kernel sizes (indicated as ‘Inception’).



# Thank You