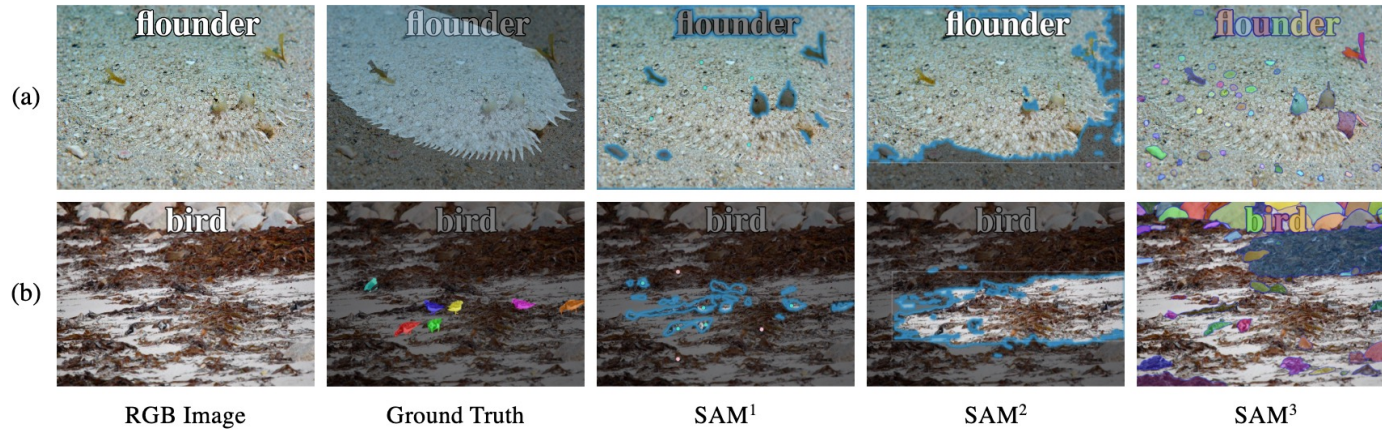


Convolution Meets LoRA: Parameter Efficient Finetuning for Segment Anything Model

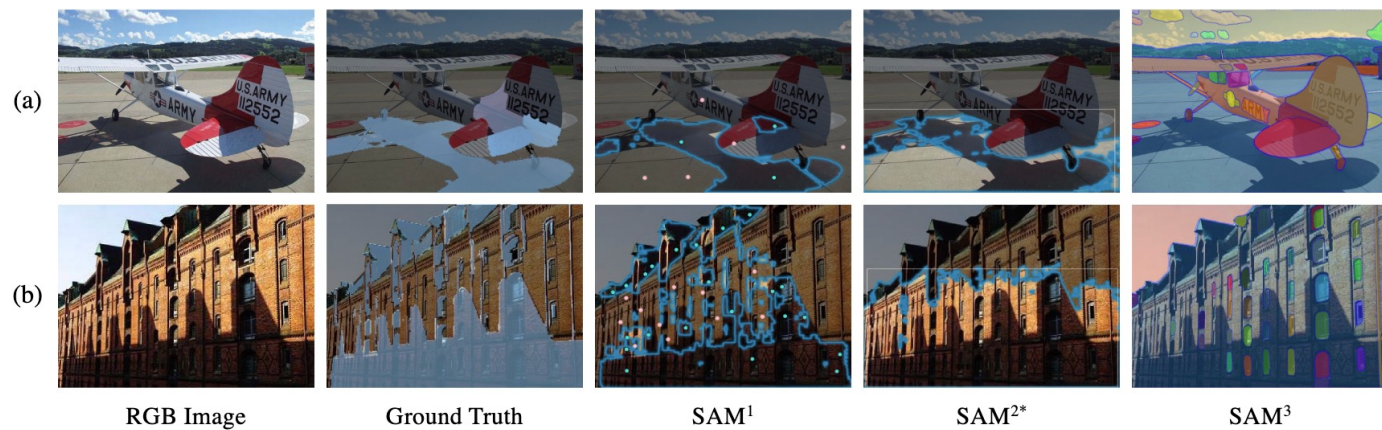
Zihan Zhong¹, Zhiqiang Tang², Tong He², Haoyang Fang², Chun Yuan¹

¹Tsinghua University, ²Amazon Web Services

Background



Segment Anything Model (SAM) may not perform well on many real-world segmentation tasks.



Background

Some existing works that fine-tuning SAM have failed to either analyze or address certain limitations inherent in SAM:

- 1) SAM's image encoder is a plain ViT, which is known to **lack of vision-specific inductive biases**.
- 2) SAM's low-level mask prediction pretraining **hinders its ability to capture high-level image semantic information**.

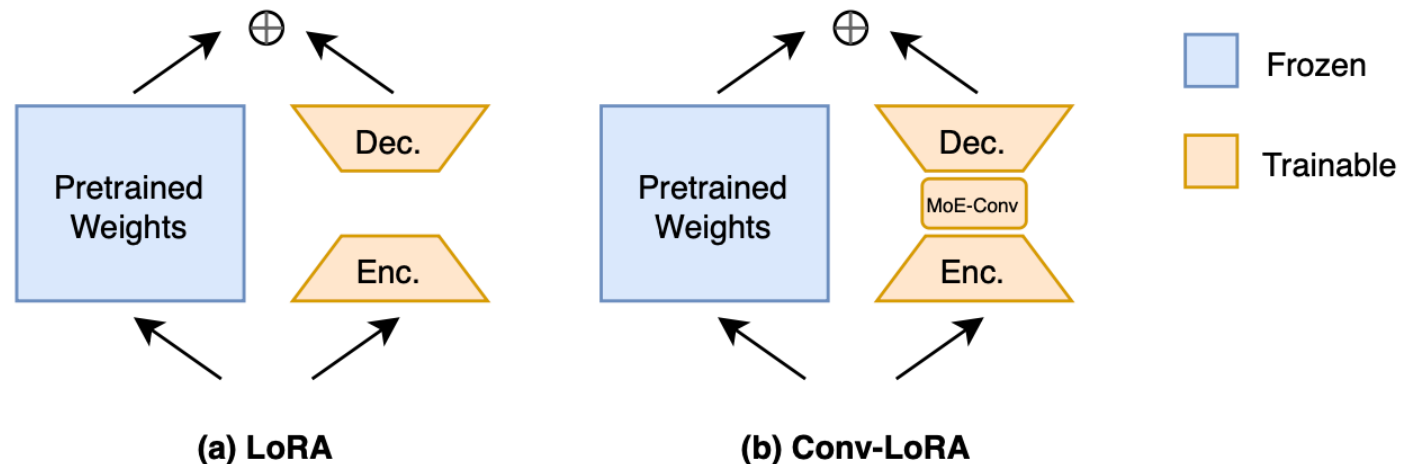
We perform linear probing for the ViT-B pretrained with Masked Auto-Encoder (MAE) and SAM on ImageNet-1K.

Method	Acc.
MAE	67.7%
SAM	54.2%

To tackle the above limitations and still retain SAM's valuable segmentation knowledge acquired during pretraining, we use **parameter efficient finetuning (PEFT)**.

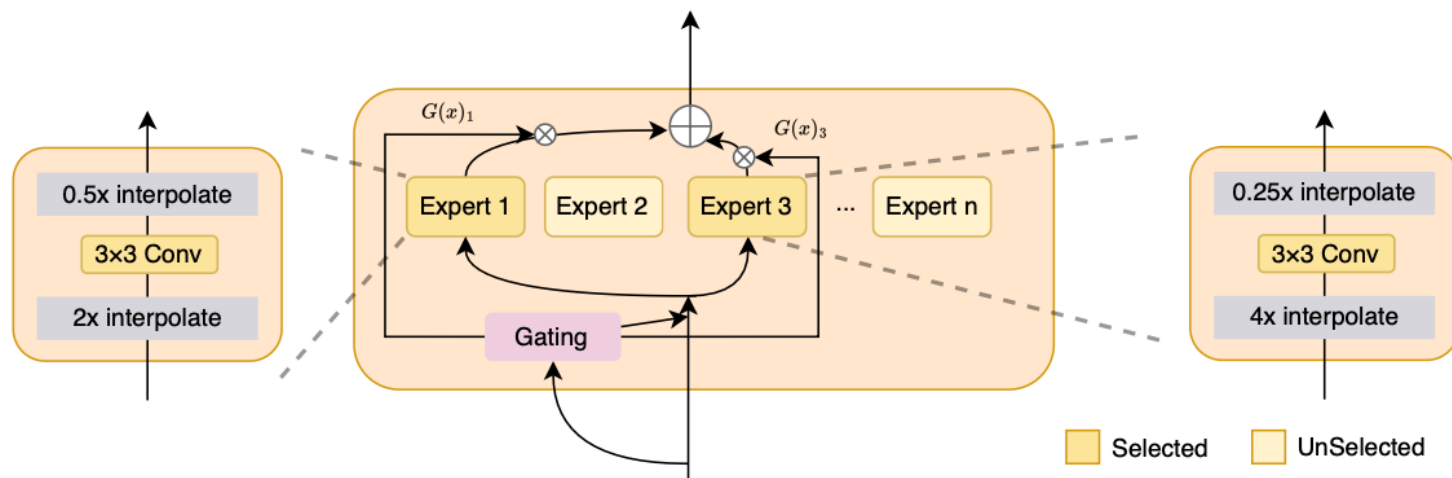
New PEFT Method Conv-LoRA

- **LoRA** surpasses widely-used visual prompt tuning (VPT), particularly in the multi-class semantic segmentation tasks.
- **Convolution operations:** inject image-related local prior.
- **Mixture-of-Experts (MoE):** dynamically select appropriate scale(s) of image features.



MoE-Conv Module

- The combination of convolution and MoE: inject the local prior into the appropriate scale(s) of image features.



Gating: dynamic expert selection.

Experts: each expert E_i specializes in one unique feature scale s_i .

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

Experiments

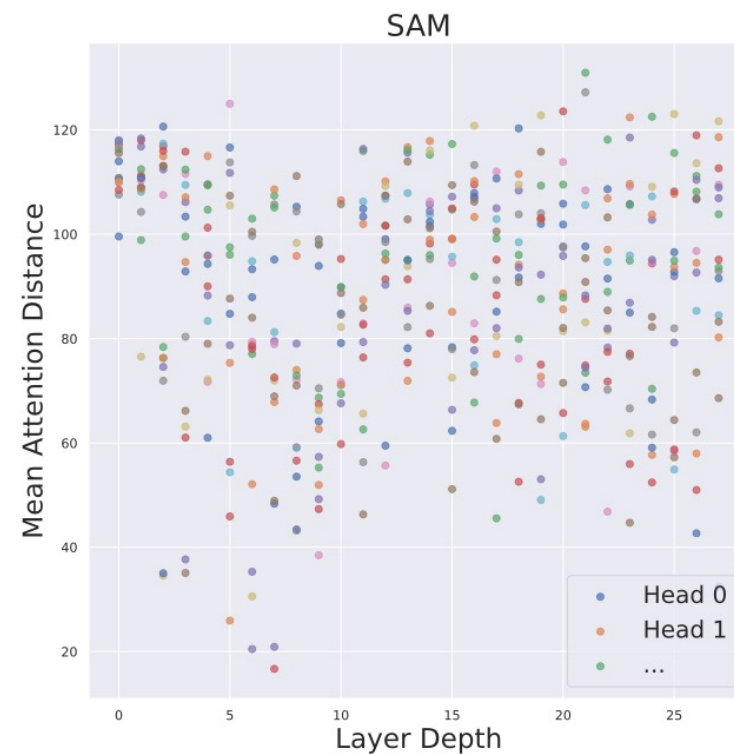
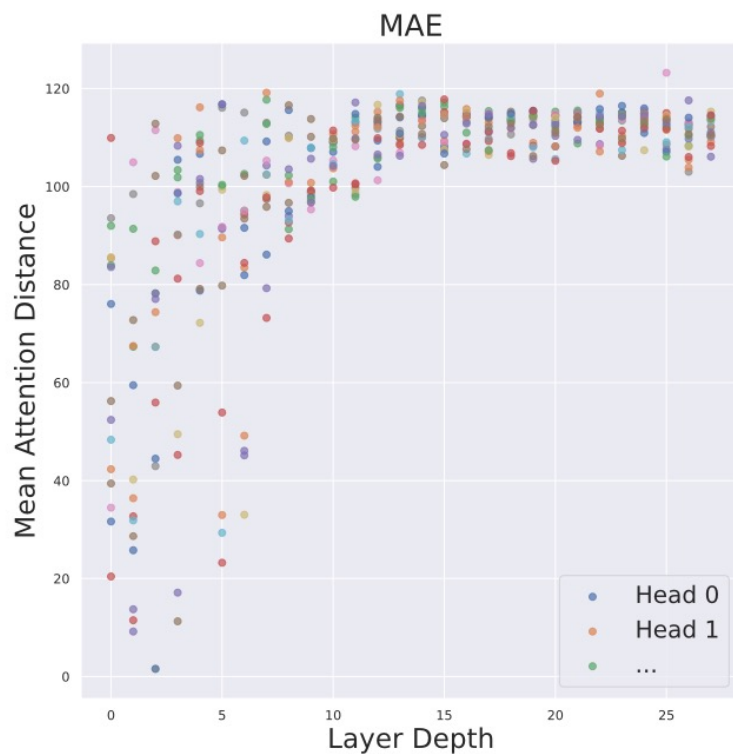
- Different Domains (Medical, Natural Images, Agriculture, Remote Sensing):
- Baselines: fine-tune SAM's decoder only, SAM trained from scratch, other PEFT methods, domain-specific models.

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
Domain Specific	* / 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
SAM trained from scratch	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 0.98	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 0.60	91.4 \pm 0.68	82.1 \pm 0.75	2.70 \pm 0.055	73.6 \pm 0.26	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 0.02	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 1.06	76.1 \pm 0.45	84.6 \pm 0.37	85.6 \pm 0.26	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 0.046	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%	92.0 \pm 0.15	94.7 \pm 0.16	91.3 \pm 0.69	94.0 \pm 0.78	77.6 \pm 0.57	85.7 \pm 0.36	88.3 \pm 0.40	92.4 \pm 0.31	84.0 \pm 0.34	2.54 \pm 0.081	74.5 \pm 0.39	84.3 \pm 0.34	62.6 \pm 0.36	76.8 \pm 0.27

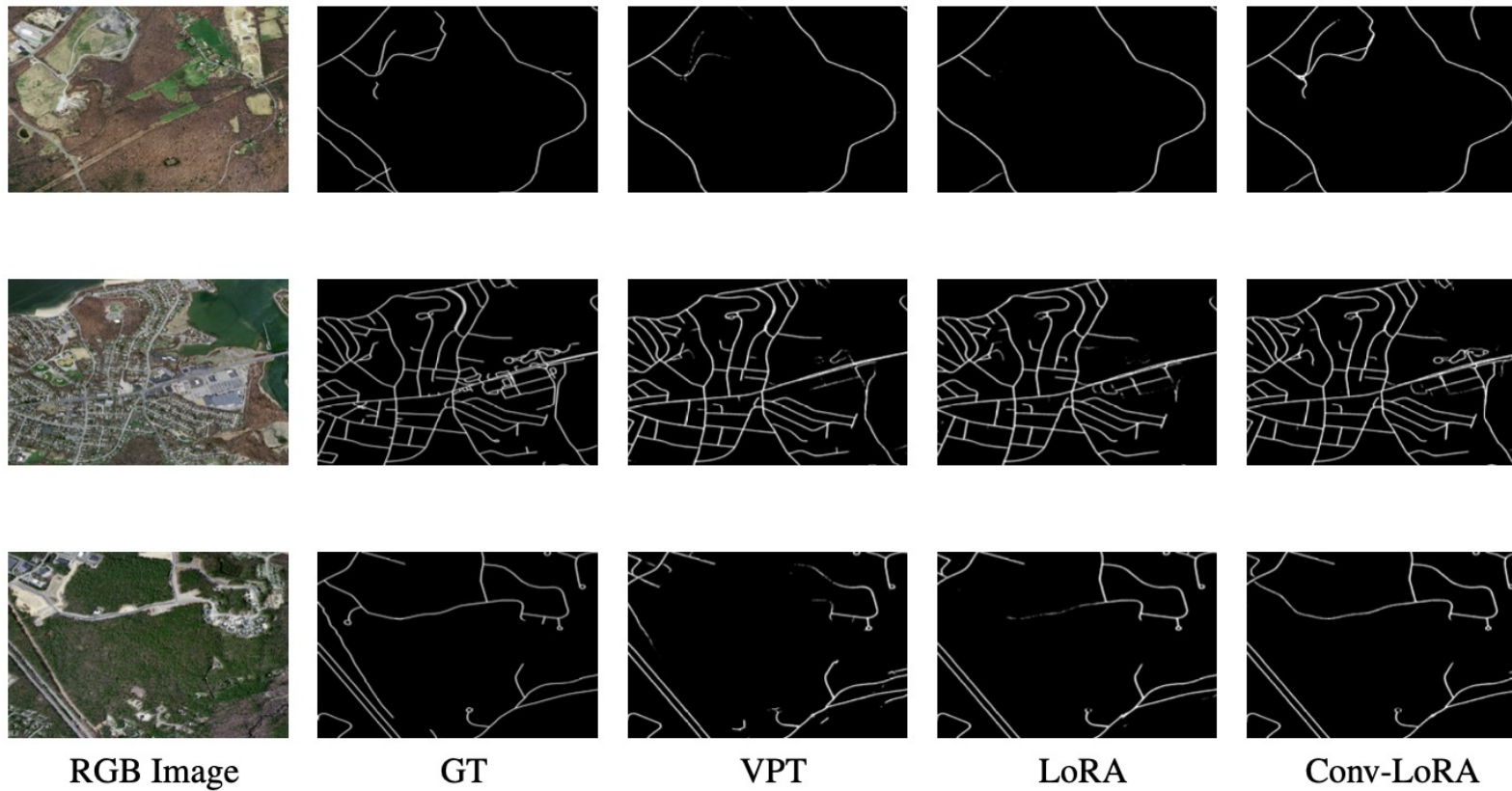
Method	# Params (M) / Ratio (%)	Acc \uparrow	mIoU \uparrow	Category IoU \uparrow											
				bg	shelf	jar	freezer	window	door	eyeglass	cup	wall	bowl	bottle	box
TransLab	42.19 / 100%	92.67	69.00	93.90	54.36	64.48	<u>65.14</u>	54.58	57.72	79.85	81.61	72.82	69.63	77.50	56.43
Trans2Seg	56.20 / 100%	94.14	<u>72.15</u>	95.35	<u>53.43</u>	<u>67.82</u>	<u>64.20</u>	<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	64.35	53.44	64.08	87.28	85.28	80.43	63.67	77.97	49.56
Conv-LoRA	4.02 / 0.63%	95.07	67.09	97.66	50.51	58.44	51.70	55.69	65.22	85.23	84.84	80.97	72.84	79.83	52.73

Experiments

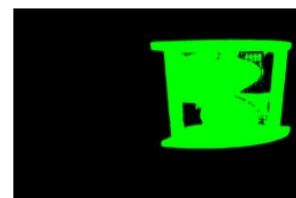
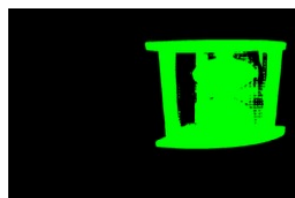
- SAM's local prior assumption: Through supervised segmentation pre-training on a vast dataset, SAM has honed a robust capability to discern and capture local features within images.



Visualization



Visualization



RGB Image

GT

VPT

LoRA

Conv-LoRA

Thank you!