PG-SAM: Prior-Guided SAM with Medical for Multi-organ Segmentation

Yiheng Zhong^{†1,2,3}, Zihong Luo^{†1,2,3}, Chengzhi Liu^{2,3}, Feilong Tang^{1,4}, Zelin Peng⁵, Ming Hu⁴, Yingzhen Hu², Jionglong Su², Zongyuan Ge^{⋈4}, and Imran Razzak^{⋈1}

- $^{1}\,$ Mohamed bin Zayed University of AI, Abu Dhabi, UAE
 - ² Xi'an Jiaotong-Liverpool University, Suzhou, China
 - University of Liverpool, Liverpool, United Kingdom Monash University, Melbourne, Australia
 - Shanghai Jiao Tong University, Shanghai, China imran.razzak@mbzuai.ac.ae

Abstract. Segment Anything Model (SAM) demonstrates powerful zeroshot capabilities; however, its accuracy and robustness significantly decrease when applied to medical image segmentation. Existing methods address this issue through modality fusion, integrating textual and image information to provide more detailed priors. In this study, we argue that the granularity of text and the domain gap affect the accuracy of the priors. Furthermore, the discrepancy between high-level abstract semantics and pixel-level boundary details in images can introduce noise into the fusion process. To address this, we propose Prior-Guided SAM (PG-**SAM**), which employs a fine-grained modality prior aligner to leverage specialized medical knowledge for better modality alignment. The core of our method lies in efficiently addressing the domain gap with fine-grained text from a medical LLM. Meanwhile, it also enhances the priors' quality after modality alignment, ensuring more accurate segmentation. In addition, our decoder enhances the model's expressive capabilities through multi-level feature fusion and iterative mask optimizer operations, supporting unprompted learning. We also propose a unified pipeline that effectively supplies high-quality semantic information to SAM. Extensive experiments on the Synapse dataset demonstrate that the proposed PG-SAM achieves state-of-the-art performance. Our code is released at https://github.com/logan-0623/PG-SAM.

Keywords: SAM \cdot Prompt-free Multi-organ Segmentation \cdot LLM

1 Introduction

Multi-organ segmentation is a core task in medical image analysis, aiming to accurately separate multiple organs. Segment Anything Model (SAM) [21] demonstrating its broad application potential [13,16,19,24,25,31,33]. The success of

[†]Equal contribution.

 $[\]square$ Corresponding author.

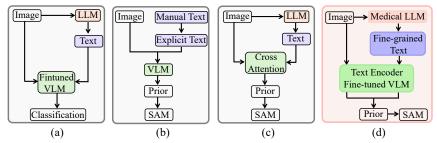


Fig. 1: Comparison of PG-SAM with other methods: (a) Issues with text granularity. (b) Fine-grained explicit text relies on manual verification and suffers from alignment problems. (c) Faces text granularity issues and lacks the zero-shot capabilities of VLM. (d) Our pipeline improves modality alignment through VLM fine-tuning, providing fine-grained text and the most comprehensive improves.

SAM relies on precise prompts. However, traditional SAM methods are time-consuming, rely on domain expertise, and are prone to human error [12,23,32,37]. To address these challenges, recent research has focused on prompt-free methods to offer simpler and more efficient segmentation solutions [8,14,34,36,30]. These methods leverage prior information to aid the decoder in better segmentation.

Inspired by multimodal learning, recent methods exploit textual information to generate priors for enhancing segmentation [1,5,18,23]. However, these approaches overlook the fact that the semantic representations derived from visual–language models (VLMs) are largely abstract and non-pixel-level, which can introduce noise. Moreover, the granularity of text descriptions influences the quality of the priors; coarse descriptions result in poor alignment with image features, thereby undermining segmentation accuracy [23]. In this work, we propose that narrowing the gap between semantic information and pixel-level boundary information can improve segmentation and mitigate noise.

Text-visual alignment methods have shown potential in guiding SAM. Prompt learning-based methods [11] generate visual descriptions by aligning with CLIP [27], providing valuable text-visual information. However, these methods suffer from text granularity limitations, which compromise modality alignment. Furthermore, the absence of a dedicated segmentation stage makes capturing fine-grained image details challenging, as illustrated in Fig. 1 (a). In contrast, TP-DRSeg [23] directly aligns explicit text with CLIP and generates priors to assist SAM, offering finer-grained information. However, it relies on ophthalmology texts verified by experts and faces alignment issues due to the VLM's reliance on natural-image training, as shown in Fig. 1 (b). Meanwhile, SEG-SAM [18] employs LLMs to provide text information, leveraging cross-attention to calculate similarity and more efficiently introduce semantic information into SAM, as depicted in Fig. 1 (c). However, it still faces granularity issues and does not fully exploit the zero-shot capability of VLM, potentially affecting its generalization.

To this end, we propose Prior-Guided SAM (PG-SAM), an efficient pipeline that provides domain-adapted, fine-grained priors and mitigates domain gap issues, as shown in Fig. 1 (d). Specifically, we introduce a fine-grained modality

prior aligner that leverages medical LLMs to merge the fluency of large-scale models with the domain-specific expertise of medical professionals, thereby excelling in complex medical scenarios and delivering more detailed, specific semantic information to CLIP. Furthermore, we fine-tune CLIP with Low-Rank Adaptation (LoRA) [15] for the medical domain, providing more accurate semantic priors. While these priors complement the embeddings of SAM, image features of CLIP focus on abstract semantics, lacking pixel-level details, which may cause simple fusion methods to blur SAM's embeddings and hinder the retention of fine-grained features. To address this, we design a novel decoder that enhances feature extraction through multi-level fusion, reducing detail loss from noise and promoting knowledge sharing between the CLIP and SAM. Finally, by leveraging an iterative mask optimizer, we dynamically fine-tune the mask weights for each category, enhancing feature expression and enabling better discrimination of small organ details. Extensive experiments on the Synapse dataset demonstrate that the proposed PG-SAM achieves state-of-the-art performance.

Overall, our contributions are threefold: (1) We propose a fine-grained modality prior aligner that combines high-level semantics and visual information to generate high-quality priors for all categories; (2) We introduce a novel decoder, which improves mask quality through multi-level feature fusion and a mask fine-tuner; (3) We provide a unified pipeline that simplifies the process while enriching prompt-free methods with medical knowledge. Experimental results demonstrate that our method improves multi-organ segmentation, outperforming state-of-theart performance on the Synapse dataset [22].

2 Methodology

2.1 Overview

Fig. 2 illustrates the overview of our method, consisting of three coordinated key components: a Fine-Grained Modality Prior Aligner, as described in Section 2.2, a Multi-level Feature Fusion, detailed in Section 2.3, and a Iterative Mask Optimizer, explained in Section 2.4. PG-SAM first generates fine-grained text descriptions for each image, combining them to generate semantic priors, referred to as the Semantic Guide Matrix, to assist the decoding process. Then the Multi-level Feature Fusion module facilitate knowledge sharing between the text-guided explicit prior and multi-level visual features. Finally, in the iterative mask optimizer, candidate masks are provided for each category via mask tokens, and a mask refiner optimizes these segmentation details.

2.2 Fine-Grained Modality Prior Aligner

The aligner employs a Medical-LLM to generate anatomically precise text prompts, whose clinical specificity enhances semantic guidance as evidenced by the sharpened heatmap patterns in Fig. 3. Then, the aligner bridges medical imaging and text domains through four key operations: First, a LoRA-tuned SAM encoder extracts multi-scale visual features $\mathbf{F}_{\text{sam}} \in \mathbb{R}^{B \times C \times H \times W}$, while a CLIP encoder processes Medical-LLM enhanced text prompts into embeddings $\mathbf{F}_{\text{text}} \in \mathbb{R}^{B \times d_{\text{text}}}$, where d_{text} represents the dimensionality of the text features and B denoted as

Fig. 2: Overview of PG-SAM. (a) Illustrates the process by which the fine-grained modality prior aligner generates the Semantic Guide Matrix G; (b) For multilevel feature fusion, G is integrated with the feature map after multi-level sampling to preserve more detailed features; (c) It outlines the iterative mask optimizer, which dynamically learns convolution kernel parameters via a Hypernetwork and refines the final mask using a dedicated refiner.

the batch size. This process ensures the accurate capture of fine-grained semantic information, providing a solid foundation for subsequent cross-modal alignment. Next, we compute dynamic similarity weights:

$$\mathbf{W}_{s} = \sigma(\mathcal{P}\left(\frac{\mathbf{F}_{\text{img}}^{\top} \mathbf{F}_{\text{text}}}{\|\mathbf{F}_{\text{img}}\| \|\mathbf{F}_{\text{text}}\|}\right)), \tag{1}$$

where \mathcal{P} denotes a learnable projection that maps CLIP's global similarity to spatial-wise weights and $\mathbf{W}_s \in \mathbb{R}^{B \times 1 \times L}$. This design explicitly quantifies cross-modal semantic alignment through cosine similarity measurement. Then we construct the spatial attention matrix $\mathbf{A} \in \mathbb{R}^{B \times L \times L}$. It captures inter-pixel relationships, enhanced through layer-normalized dot-product attention:

$$\mathbf{A} = \operatorname{softmax} \left(\frac{\mathbf{F}_{\operatorname{sam}}^{\operatorname{norm}} \cdot (\mathbf{F}_{\operatorname{sam}}^{\operatorname{norm}})^{\top}}{\sqrt{C}} \right), \tag{2}$$

where C represents the number of channels in the feature map.

Finally, the final guidance matrix ${\bf G}$ combines attention-refined features with similarity weights through channel-wise broadcasting and dual-level normalization:

$$\mathbf{G} = \Gamma_{\text{spatial}}(\Gamma_{\text{channel}}((\mathbf{F}_{\text{sam}} + \mathbf{A}\mathbf{F}_{\text{sam}}) \odot \mathbf{W}_s)), \tag{3}$$

where Γ denotes layer normalization operators, and \odot represents element-wise multiplication with broadcasting.

2.3 Multi-level Feature Fusion

Existing approaches that directly predict masks on low-resolution feature maps often lead to blurry boundaries [35], while simple bilinear upsampling loses crucial high-frequency details [29]. To address these issues, we propose a multi-level fusion module based on learnable feature reorganization.

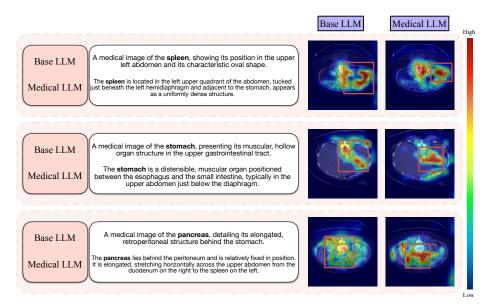


Fig. 3: Comparison of textual prompts and corresponding heatmaps generated by the Base LLM (left) and the Medical LLM (right) for anatomical images of the *spleen*, *stomach*, and *pancreas*. The Medical LLM provides clinically precise descriptions, yielding more focused and detailed semantic guidance, as demonstrated by the sharper heatmap regions.

We implement cross-scale feature fusion with a pyramid upsampling architecture, as shown in Fig. 1, and employ bilinear spatial expansion through a stride-2 transposed convolution. The two-stage upsampling process is formulated as:

$$\mathbf{F}_{\mathrm{up}}^{(t)} = \sigma\left(\mathrm{LN}\left(\mathrm{DeConv}_{2\times 2}(\mathbf{F}_{\mathrm{trans}})\right)\right), \quad t \in \{1, 2\},\tag{4}$$

where $\mathbf{F}_{\text{trans}}$ is the Transformer's output feature map. This hierarchical design with $2\times$ resolution increments per stage combined with LN-GELU modules effectively regulates gradient flow, significantly suppressing checkerboard artifacts compared to single-step upsampling [26]. The channel compression ratio of 4:2:1 is intentionally adopted to enable the stepwise recovery of high-frequency details, while maintaining computational efficiency throughout the process.

To further enhance spatial awareness, we employ deformable convolution to precisely align the guidance matrix $G \in \mathbb{R}^{C \times H \times W}$ with the upsampled features along the spatial dimensions, thereby effectively capturing complementary cues such as edges and contours that are often missed in coarse features [10]. Subsequently, a 1×1 convolution is applied to compress the channel dimension, enabling the efficient fusion of the guidance information with the upsampled features, which ultimately achieves integrated cross-modal feature enhancement:

$$F_{fusion} = \phi \left(F_{up}^{(2)} \right) + \psi \left(\text{Align}(G; \theta) \right), \tag{5}$$

where $\phi(\cdot)$ denotes a 1×1 convolution used for channel reduction, $\psi(\cdot)$ represents an affine transformation applied to the guidance matrix, and θ comprises the learnable deformation parameters.

2.4 Iterative Mask Optimizer

To address coarse edges in initial mask predictions [35], we propose an iterative mask optimizer with two core components:

Instance-Adaptive Kernel Generation. To balance general feature extraction with instance-specific adaptation, we design a hypernetwork that generates dynamic convolution parameters. For an instance i with a mask encoding $m_i \in \mathbb{R}^C$, a MLP generates dynamic convolution kernel parameters Ω_i through:

$$\Omega_i = \text{MLP}(m_i) \odot \mathcal{W}_{base}, \text{ where } \Omega_i \in \mathbb{R}^{C_{in} \times C_{out} \times K \times K}.$$
 (6)

In this equation $W_{base} \in \mathbb{R}^{C_{in} \times C_{out} \times K \times K}$ represents shared base kernels, and \odot represents the channel-wise Hadamard product [20]. This design achieves adaptability via two key aspects: (1) the base kernel W_{base} extracts universal features across diverse instances; and (2) the dynamic weights $\text{MLP}(m_i)$ encode instance-specific geometric information, modulating channel responses through a gating mechanism to effectively accommodate various object morphologies.

Progressive Residual Refinement. We implement mask optimization through iterative residual corrections. The operations of t-th iteration:

$$\Delta M_t = \sigma \left(\operatorname{Conv}_{3\times 3}^{\Omega_t} \left([M_t \oplus F_{fusion}] \right) \right), \quad M_{t+1} = \operatorname{Clip}(M_t + \lambda \Delta M_t),$$
 (7)

where M represents mask, \oplus denotes channel concatenation, λ is a learnable step-size coefficient, and $Clip(\cdot)$ constrains the output values to the range [0,1]. **Training Objective.** To train our segmentation model, the overall training objective adopts the combination of cross-entropy and Dice similarity:

$$\mathcal{L}_{\text{loss}} = \sum_{r \in l, h} \left[(1 - \lambda) \mathcal{L}_{\text{CE}}^{r} + \lambda \mathcal{L}_{\text{Dice}}^{r} \right], \tag{8}$$

where $r \in \{l, h\}$ denotes low/high-resolution paths (56×56 and 224×224) and a hyperparameters $\lambda = 0.8$ controls their balance.

3 Experiment

3.1 Experimental Setup

Dataset. We evaluate on the MICCAI 2015 Synapse Multi-Organ CT datase [22] containing 3,779 contrast-enhanced abdominal CT slices (2,212 training). Following SAMed [36] and H-SAM [9], we use 18/12 cases for training/test with 224×224 resolution slices. Evaluation covers eight organs: aorta, gallbladder, spleen, kidneys, liver, pancreas, and stomach.

Implementation details. We implement training on an RTX 4090 GPU with H-SAM-compatible augmentations. The maximum number of training epochs is set to 300, and the AdamW optimizer is used with β_1 , β_2 , and weight decay set to 0.9, 0.999, and 0.1, respectively. Additionally, we follow the same LoRA configuration as SAMed, where the rank of LoRA is set to 4.

Table 1: Comparison with state-of-the-art models on Synapse multi-organ CT dataset in both few-shot and fully-supervised settings. Greyed values represent our results, while bold values indicate outperformance of SOTA models. mDice: the Mean Dice coefficient; HD95: the 95th percentile of the Hausdorff Distance.

	Method	Spleen	Kidney(R)	Kidney(L)	Gallbladder	Liver	Stomach	Aorta	Pancreas	mDice†	HD95↓
10%	AutoSAM [17]	68.80	77.44	76.53	24.87	88.06	52.70	75.19	34.58	55.69	31.67
	SAM Adapter [7]	72.42	68.38	66.77	22.38	89.69	53.15	66.74	26.76	58.28	54.42
	SAMed [36]	87.32	80.10	82.75	70.24	93.37	73.62	86.99	67.64	80.26	28.89
	H-SAM [9]	90.87	83.89	81.99	61.59	93.69	76.07	83.26	50.92	77.79	18.03
	PG-SAM	88.43	82.06	82.23↑	53.75	92.27	78.80↑	82.04	46.43	75.75	$12.35 \uparrow$
	TransUnet [6]	87.23	63.13	81.87	77.02	94.08	55.86	85.08	75.62	77.48	31.69
Fully - Supervised	SwinUnet [4]	85.47	66.53	83.28	79.61	94.29	56.58	90.66	76.60	79.13	21.55
	TransDeepLab [3]	86.04	69.16	84.08	79.88	93.53	61.19	89.00	78.40	80.16	21.25
	DAE-Former [2]	88.96	72.30	86.08	80.88	94.98	65.12	91.94	79.19	82.43	17.46
	MERIT [28]	92.01	84.85	87.79	74.40	95.26	85.38	87.71	71.81	84.90	13.22
	AutoSAM [17]	80.54	80.02	79.66	41.37	89.24	61.14	82.56	44.22	62.08	27.56
	SAM Adapter [7]	83.68	79.00	79.02	57.49	92.68	69.48	77.93	43.07	72.80	33.08
	SAMed [36]	87.33	80.10	82.75	70.24	93.37	73.62	86.99	67.64	80.26	28.89
	H-SAM [9]	92.34	85.99	87.71	69.65	95.20	86.27	87.53	72.53	84.65	7.29
	PG-SAM	93.12	84.57	87.93	73.26	95.40	86.62	87.87	71.49	84.79↑	7.61
_											

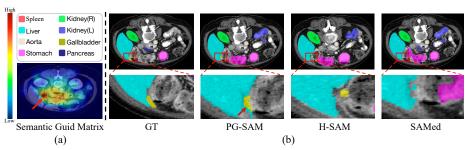


Fig. 4: (a) Shows one of the focused areas in our semantic guide matrix; (b) Displays the visualization of segmentation results from various methods on the Synapse dataset especially focus on gallbladder.

3.2 Comparisons with State-of-the-art Methods

As shown in Table 1, PG-SAM shows substantial improvements across both few-shot and fully supervised scenarios. Under the 10% annotation setting, our method exceeds state-of-the-art performance in left kidney segmentation ($\uparrow 0.24\%$) and stomach segmentation ($\uparrow 2.73\%$), while reducing boundary localization errors to HD95=12.35($\downarrow 5.68\%$) compared to the best-performing baseline, demonstrating superior boundary-aware segmentation capabilities. In the fully supervised setting, PG-SAM attains the highest mean Dice coefficient ($\uparrow 0.14\%$) among prompt-free SAM variants, with particularly notable improvements in challenging anatomical structures: spleen ($\uparrow 0.78\%$), left kidney ($\uparrow 0.22\%$), and gallbladder ($\uparrow 3.61\%$). Compared to conventional fully supervised methods, PG-SAM achieves dual improvements in both segmentation accuracy and boundary precision: 1) Surpasses TransUNet ($\uparrow 7.31\%$ Dice) and SwinUNet ($\uparrow 5.66\%$ Dice) in overall segmentation quality, 2) Reduces boundary errors by HD95=7.61 to MERIT ($\downarrow 5.61\%$), while maintaining comparable Dice performance ($\downarrow 0.11\%$).

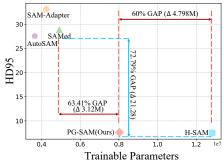


Fig. 5: Comparing HD95 scores against trainable parameters for different SAMs, with GAP defined as percentage change.

Table 2: Ablation study on the effectiveness of key components in PG-SAM: Fine-Grained Modality Prior Aligner (FGMPA), Multi-level Feature Fusion (MLFF), and Iterative Mask Optimizer (IMO), in terms of mean Dice (mDice) (%).

	FGMPA	MLFF	IMO	mDice (%)
Ι				72.80%
Π	✓			77.28%
III	✓	✓		80.10%
${\rm IV}$	✓	✓	✓	84.79%

This demonstrates our method's unique strength in achieving precise boundary delineation without compromising region-wise segmentation accuracy.

Additionally, PG-SAM achieves remarkable performance efficiency, as evidenced by our experimental results, as demonstrated in Fig. 5. When comparing Hausdorff Distance 95 (HD95) scores against the number of trainable parameters, PG-SAM demonstrates superior efficiency by maintaining competitive or superior segmentation accuracy while utilizing significantly fewer parameters than existing state-of-the-art approaches. Despite operating at lower resolution $(224 \times 224 \text{ vs. } 512 \times 512)$ in few-shot settings which may marginally affect fine detail capture in Medical image, PG-SAM still remains competitive overall.

3.3 Qualitative Results

In this example, we select the **gallbladder** as the region of interest for semantic guidance. Fig. 4 (a) shows a heatmap that highlights the key focus areas of the semantic guide matrix, while Fig. 4 (b) presents the corresponding segmentation results. The prior information effectively guides the localization of the gallbladder: SAMed fails to localize it, H-SAM segments it inaccurately, whereas PG-SAM both locates and accurately segments the gallbladder.

3.4 Ablation Study

To validate the efficacy of the three core components in PG-SAM: FGMPA, MLFF, and IMO, we perform ablation studies on the Synapse dataset, as shown in Table 2. Starting from Experiment I (Baseline) with a Mean Dice of 72.80%, adding FGMPA in Experiment II improves cross-modal alignment, leading to a +4.48% gain. Further integrating MLFF in Experiment III enhances feature fusion, increasing performance by +2.82%. Finally, incorporating IMO in Experiment IV refines segmentation masks through iterative optimization, achieving the highest Mean Dice of 84.79%, with a final improvement of +4.69% over Experiment III. These results demonstrate the synergistic advantage of combining all three components, each contributing to the overall segmentation accuracy.

4 Conclusion

In this study, we address the limitations of SAM in medical image segmentation, where domain gaps and insufficient textual priors lead to performance degradation. To this end, our proposed PG-SAM integrates medical LLMs to enhance segmentation accuracy. It introduces three key innovations: (1) a fine-grained modality prior aligner for precise anatomical priors, (2) a multi-level feature fusion module that seamlessly integrates global semantic context with local structural details, and (3) an iterative mask optimizer that progressively refines boundary precision. Comprehensive experiments on the Synapse dataset demonstrate that PG-SAM exceeds state-of-the-art performance, enhancing multi-organ segmentation accuracy, especially for complex organs.

References

- Aleem, S., Wang, F., Maniparambil, M., Arazo, E., Dietlmeier, J., Curran, K., Connor, N.E., Little, S.: Test-time adaptation with salip: A cascade of sam and clip for zero-shot medical image segmentation. In: CVPR. pp. 5184–5193 (2024)
- Azad, R., Arimond, R., Aghdam, E.K., Kazerouni, A., Merhof, D.: Dae-former: Dual attention-guided efficient transformer for medical image segmentation. In: International Workshop on PRedictive Intelligence In MEdicine. pp. 83–95. Springer (2023)
- Azad, R., Heidari, M., Shariatnia, M., Aghdam, E.K., Karimijafarbigloo, S., Adeli, E., Merhof, D.: Transdeeplab: Convolution-free transformer-based deeplab v3+ for medical image segmentation. In: International Workshop on PRedictive Intelligence In MEdicine. pp. 91–102. Springer (2022)
- 4. Cao, H., Wang, Y., Chen, J., Jiang, D., Zhang, X., Tian, Q., Wang, M.: Swin-unet: Unet-like pure transformer for medical image segmentation. In: European conference on computer vision. pp. 205–218. Springer (2022)
- 5. Chen, H., Xu, Y., Xu, Y., Zhang, Y., Cui, L.: Test-time medical image segmentation using clip-guided sam adaptation. In: 2024 IEEE International Conference on Bioinformatics and Biomedicine (BIBM). pp. 1866–1873. IEEE (2024)
- Chen, J., Lu, Y., Yu, Q., Luo, X., Adeli, E., Wang, Y., Lu, L., Yuille, A.L., Zhou, Y.: Transunet: Transformers make strong encoders for medical image segmentation. arXiv preprint arXiv:2102.04306 (2021)
- Chen, T., Zhu, L., Ding, C., Cao, R., Wang, Y., Li, Z., Sun, L., Mao, P., Zang, Y.: Sam fails to segment anything?—sam-adapter: Adapting sam in underperformed scenes: Camouflage, shadow, medical image segmentation, and more. arXiv preprint arXiv:2304.09148 (2023)
- 8. Chen, Z., Xu, Q., Liu, X., Yuan, Y.: Un-sam: Universal prompt-free segmentation for generalized nuclei images. arXiv preprint arXiv:2402.16663 (2024)
- 9. Cheng, Z., Wei, Q., Zhu, H., Wang, Y., Qu, L., Shao, W., Zhou, Y.: Unleashing the potential of sam for medical adaptation via hierarchical decoding. In: CVPR. pp. 3511–3522 (2024)
- Dai, J., Qi, H., Xiong, Y., Li, Y., Zhang, G., Hu, H., Wei, Y.: Deformable convolutional networks. In: Proceedings of the IEEE international conference on computer vision. pp. 764–773 (2017)

- Fang, X., Lin, Y., Zhang, D., Cheng, K.T., Chen, H.: Aligning medical images with general knowledge from large language models. In: International Conference on Medical Image Computing and Computer-Assisted Intervention. pp. 57–67. Springer (2024)
- 12. Feng, W., Zhu, L., Yu, L.: Cheap lunch for medical image segmentation by finetuning sam on few exemplars. In: International MICCAI Brainlesion Workshop. pp. 13–22. Springer (2023)
- 13. He, S., Bao, R., Li, J., Stout, J., Bjornerud, A., Grant, P.E., Ou, Y.: Computer-vision benchmark segment-anything model (sam) in medical images: Accuracy in 12 datasets. arXiv preprint arXiv:2304.09324 (2023)
- 14. He, X., Hu, Y., Zhou, Z., Jarraya, M., Liu, F.: Few-shot adaptation of training-free foundation model for 3d medical image segmentation. arXiv preprint arXiv:2501.09138 (2025)
- 15. Hu, E.J., Shen, Y., Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S., Wang, L., Chen, W., et al.: Lora: Low-rank adaptation of large language models. ICLR 1(2), 3 (2022)
- 16. Hu, M., Li, Y., Yang, X.: Skinsam: Empowering skin cancer segmentation with segment anything model. arXiv preprint arXiv:2304.13973 (2023)
- Hu, X., Xu, X., Shi, Y.: How to efficiently adapt large segmentation model (sam) to medical images. arXiv preprint arXiv:2306.13731 (2023)
- 18. Huang, S., Liang, H., Wang, Q., Zhong, C., Zhou, Z., Shi, M.: Seg-sam: Semantic-guided sam for unified medical image segmentation. arXiv preprint arXiv:2412.12660 (2024)
- 19. Ji, G.P., Fan, D.P., Xu, P., Cheng, M.M., Zhou, B., Van Gool, L.: Sam struggles in concealed scenes—empirical study on segment anything. arXiv preprint arXiv:2304.06022 (2023)
- Kim, J.H., On, K.W., Lim, W., Kim, J., Ha, J.W., Zhang, B.T.: Hadamard product for low-rank bilinear pooling. arXiv preprint arXiv:1610.04325 (2016)
- Kirillov, A., Mintun, E., Ravi, N., Mao, H., Rolland, C., Gustafson, L., Xiao, T., Whitehead, S., Berg, A.C., Lo, W.Y., et al.: Segment anything. In: Proceedings of the IEEE/CVF international conference on computer vision. pp. 4015–4026 (2023)
- 22. Landman, B.A., Xu, Z., Iglesias, J.S., Styner, M.A., Langerak, T., Klein, A.K.: Multi-atlas labeling beyond the cranial vault—workshop and challenge (2015). https://doi.org/10.7303/syn3193805
- 23. Li, W., Xiong, X., Xia, P., Ju, L., Ge, Z.: Tp-drseg: improving diabetic retinopathy lesion segmentation with explicit text-prompts assisted sam. In: MICCAI. pp. 743–753. Springer (2024)
- Li, Y., Hu, M., Yang, X.: Polyp-sam: Transfer sam for polyp segmentation. In: Medical Imaging 2024: Computer-Aided Diagnosis. vol. 12927, pp. 759–765. SPIE (2024)
- 25. Mazurowski, M.A., Dong, H., Gu, H., Yang, J., Konz, N., Zhang, Y.: Segment anything model for medical image analysis: an experimental study. Medical Image Analysis 89, 102918 (2023)
- Odena, A., Dumoulin, V., Olah, C.: Deconvolution and checkerboard artifacts. Distill (2016)
- Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al.: Learning transferable visual models from natural language supervision. In: International conference on machine learning. pp. 8748–8763. PmLR (2021)
- 28. Rahman, M.M., Marculescu, R.: Multi-scale hierarchical vision transformer with cascaded attention decoding for medical image segmentation. In: Medical Imaging with Deep Learning. pp. 1526–1544. PMLR (2024)

- 29. Ronneberger, O., Fischer, P., Brox, T.: U-net: Convolutional networks for biomedical image segmentation (2015)
- 30. Tang, F., Xu, Z., Hu, M., Li, W., Xia, P., Zhong, Y., Wu, H., Su, J., Ge, Z.: Neighbor does matter: Density-aware contrastive learning for medical semi-supervised segmentation. In: AAAI (2025)
- 31. Tang, F., Xu, Z., Qu, Z., Feng, W., Jiang, X., Ge, Z.: Hunting attributes: Context prototype-aware learning for weakly supervised semantic segmentation. In: CVPR (2024)
- 32. Wu, J., Ji, W., Liu, Y., Fu, H., Xu, M., Xu, Y., Jin, Y.: Medical sam adapter: Adapting segment anything model for medical image segmentation. arXiv preprint arXiv:2304.12620 (2023)
- 33. Wu, Z., Zhao, Q., Hu, M., Li, Y., Xue, H., Dang, K., Jiang, Z., Stefanidis, A., Wang, Q., Razzak, I., et al.: Mswal: 3d multi-class segmentation of whole abdominal lesions dataset. arXiv preprint arXiv:2503.13560 (2025)
- 34. Xie, B., Tang, H., Duan, B., Cai, D., Yan, Y.: Masksam: Towards auto-prompt sam with mask classification for medical image segmentation. arXiv preprint arXiv:2403.14103 (2024)
- 35. Zhang, G., Lu, X., Tan, J., Li, J., Zhang, Z., Li, Q., Hu, X.: Refinemask: Towards high-quality instance segmentation with fine-grained features (2021)
- Zhang, K., Liu, D.: Customized segment anything model for medical image segmentation. arXiv preprint arXiv:2304.13785 (2023)
- 37. Zhang, Y., Leng, T., Han, K., Xie, X.: Self-sampling meta sam: Enhancing few-shot medical image segmentation with meta-learning. arXiv preprint arXiv:2308.16466 (2023)