

MedSAM-CA: A CNN-Augmented ViT with Attention-Enhanced Multi-Scale Fusion for Medical Image Segmentation

Peiting Tian^a, Xi Chen^{a,*}, Haixia Bi^a, Fan Li^a

^a*School of Information and Communications Engineering, Xi'an Jiaotong University, 710049, Xi'an, 28 Xianning West Road, Shaanxi, China*

Abstract

Medical image segmentation plays a crucial role in clinical diagnosis and treatment planning, where accurate boundary delineation is essential for precise lesion localization, organ identification, and quantitative assessment. In recent years, deep learning-based methods have significantly advanced segmentation accuracy. However, two major challenges remain. First, the performance of these methods heavily relies on large-scale annotated datasets, which are often difficult to obtain in medical scenarios due to privacy concerns and high annotation costs. Second, clinically challenging scenarios, such as low contrast in certain imaging modalities and blurry lesion boundaries caused by malignancy, still pose obstacles to precise segmentation. To address these challenges, we propose MedSAM-CA, an architecture-level fine-tuning approach that mitigates reliance on extensive manual annotations by adapting the pretrained foundation model, Medical Segment Anything (MedSAM). MedSAM-CA introduces two key components: the Convolutional Attention-Enhanced Boundary Refinement Network (CBR-Net) and the Attention-Enhanced Feature Fusion Block (Atte-FFB). CBR-Net operates in parallel with the MedSAM encoder to recover boundary information potentially overlooked by long-range attention mechanisms, leveraging hierarchical convolutional processing. Atte-FFB, embedded in the MedSAM decoder, fuses multi-level fine-grained features from skip connections in CBR-Net with global representations upsampled within the decoder to enhance boundary delineation accuracy. Experiments on publicly available datasets

*Corresponding author: xi_chen@mail.xjtu.edu.cn

covering dermoscopy, CT, and MRI imaging modalities validate the effectiveness of MedSAM-CA. On dermoscopy dataset for melanoma segmentation, MedSAM-CA achieves 94.43% Dice with only 2% of full training data, reaching 97.25% of full-data training performance, demonstrating strong effectiveness in low-resource clinical settings.

Keywords:

medical image segmentation, MedSAM foundation model, fine-tuning, multi-scale fusion

1. Introduction

High-precision medical image segmentation techniques can support disease diagnosis, radiotherapy and chemotherapy dose planning, and treatment efficacy evaluation, serving as an indispensable technological cornerstone for advancing precision clinical medicine [1]. Precise delineation of lesion boundaries is essential in clinical decision-making, as boundary information provides critical structural cues that facilitate tasks such as confirming visceral metastases and determining therapeutic margins for melanoma. To reduce repetitive workload, researchers have developed computer-aided diagnosis (CAD) systems to support efficient and accurate lesion and organ segmentation [2].

Recent advances in deep learning have shown great promise in automating medical image segmentation [3]. However, their performance remains constrained by the need for large-scale and diverse annotated datasets. In clinical practice, such datasets are difficult to obtain due to factors such as patient privacy, high annotation costs, and the scarcity of expert-verified labels, making it challenging to meet the data requirements for training robust models. In addition, accurate automatic boundary segmentation in medical images remains a significant challenge due to intrinsic factors such as heterogeneous lesion morphology, limited imaging resolution, and modality-specific constraints.

To mitigate this limitation, large-scale pretraining models enable to learn transferable representations and reduce reliance on task-specific annotations in resource-limited settings [4]. Building on this, the Segment Anything Model (SAM), trained on 11 million images and over 1 billion segmentation masks, has demonstrated strong zero-shot transfer capabilities to unseen image distributions and tasks [5]. SAM demonstrates strong perfor-

mance on natural images. Medical Segment Anything Model (MedSAM), built upon SAM and further trained on large-scale medical datasets comprising 1,570,263 image-mask pairs [6], has emerged as one of the most promising foundation models for next-generation medical image segmentation.

SAM and its medical adaptation, MedSAM, have alleviated data dependency to some extent. However, their performance remains limited in domain-specific clinical scenarios, particularly when training data is extremely scarce or when segmenting low-contrast lesions and fine anatomical boundaries. To mitigate these shortcomings, researchers have proposed different SAM-based strategies, such as Full-Parameter Supervised Fine-Tuning (Full SFT), Parameter-Efficient Fine-Tuning (PEFT), prompt engineering, domain-specific self-supervised pre-training, and architecture-level fine-tuning.

Full SFT retrains all MedSAM weights end-to-end, enabling the model to specialize in rare or complex conditions with limited annotations; it has achieved notable success in tasks such as impacted-tooth detection [7] and mediastinal segmentation [8]. PEFT techniques, exemplified by lightweight Adapters and Low-Rank Adaptation (LoRA), insert only a small set of trainable parameters into an otherwise frozen backbone, thereby reducing computational and memory overhead. For instance, Wu et al. [9] propose an approach named Med-SA, embedding Adapter blocks into SAM for three-dimensional (3D) segmentation enhancement, whereas Li et al. [10] couple LoRA with convolutional prediction heads to achieve comparable gains without full retraining. Given the semi-automated nature of SAM-based models, prompt engineering constitutes another an important research area. One-Prompt Former extends a single annotated prompt across diverse segmentation tasks, whereas self-prompting methods learn to synthesize point-type and box-type inputs directly from encoder embeddings, enabling few-shot segmentation without human intervention [11, 12]. Domain-specific self-supervised pretraining methods, such as Swin-UNETR and MedLSAM, enhance generalizability and reduce reliance on annotations by leveraging anatomical priors through multi-task learning or implicit coordinate encoding [13, 14]. Architecture-level refinements enhance downstream segmentation performance by appending lightweight auxiliary branches or decoders. I-MedSAM augments the SAM encoder with a lightweight high-frequency adapter and an implicit-neural-representation (INR) decoder, tightening boundary contours while keeping trainable parameters to about 1.6 M [15]. Likewise, Li et al. retain the pretrained MobileSAM Vision Transformer

(ViT) as a frozen feature extractor and attach a parallel nnU-Net pathway with a curvature-aware regression head to sharpen edges [16].

Among existing fine-tuning strategies, architecture-level fine-tuning offers inherent advantages for enhancing boundary segmentation accuracy. By introducing or modifying structural components based on the original architecture, this approach extends the model’s representational capacity rather than merely improving parameter efficiency through partial tuning. Recent implementations of this strategy, including I-MedSAM and nnSAM [15, 16], aim to enhance boundary delineation by replacing the native decoder of SAM or MedSAM. However, this replacement disrupts the original encoder-decoder information flow and discards pretrained decoder parameters, potentially compromising the model’s generalization ability and limiting its performance on unseen medical data. Furthermore, these introduced decoders often require substantial amounts of training data to exceed the segmentation performance of the MedSAM.

To address this limitation, we propose MedSAM-CA, a segmentation framework designed to enhance MedSAM’s boundary sensitivity while enabling fine-tuning for domain-specific tasks with minimal training data. MedSAM-CA incorporates three key architectural components: the Convolutional Attention-Enhanced Boundary Refinement Network (CBR-Net), the Attention-Enhanced Feature Fusion Block (Atte-FFB), and an embedded Adapter module.

CBR-Net operates in parallel with the MedSAM encoder through a lightweight convolutional branch that leverages the locality and inductive bias of Convolutional Neural Networks (CNNs) to capture fine-grained details. Skip connections are employed to restore shallow-layer information, thereby mitigating boundary degradation. Furthermore, the Atte-FFB adaptively fuses local and global features across multiple scales, enhancing both segmentation accuracy and generalizability. Finally, an Adapter module embedded within the MedSAM encoder enables task-specific adaptation without modifying the pretrained parameters.

Our main contributions are summarized as follows:

We propose MedSAM-CA, a lightweight structural fine-tuning framework tailored for MedSAM, which effectively addresses the challenges of limited annotated data and low boundary accuracy. By fine-tuning only 7M parameters, MedSAM-CA achieves 94.43% Dice with only 2% of training data on melanoma segmentation, reaching 97.25% of full-data performance.

We design CBR-Net, a parallel CNN branch that enhances local feature

extraction and boundary delineation while preserving most of the pretrained ViT encoder. The skip connections help retain low-level spatial information, leading to significantly improved boundary precision.

We design Atte-FFB, a feature fusion block that adaptively integrates multi-scale local and global representations, thereby improving segmentation accuracy and enhancing model generalization across imaging modalities.

2. Related works

2.1. Medical image segmentation

Before the advent of machine learning, traditional segmentation algorithms played an essential role in automated medical image analysis. These included edge-based methods relying on closed contours [17], statistical approaches grounded in shape and appearance modeling [18], and region-based techniques guided by predefined rules [19]. Physics-driven models such as active contour methods [20] also contributed, particularly to 3D segmentation tasks in Computed Tomography (CT) imaging.

The emergence of fully convolutional networks (FCNs), pioneered by Long et al. [21], marked a turning point in deep learning-based segmentation. CNNs proved highly effective in capturing complex spatial hierarchies and subtle boundary information, making them particularly suitable for medical image segmentation. Among CNN variants, U-Net [22] became a foundational architecture due to its encoder–decoder design with skip connections, which preserve spatial detail. Building on U-Net, numerous extensions, including ResU-Net [23], U-Net++ [24], V-Net [25], 3D U-Net [26], and DeepLabV3+ [27] have further advanced segmentation performance across various modalities. For example, Hermes-R [28] combines a ResU-Net backbone with oracle-guided context-prior learning, achieving accurate segmentation of abdominal organs.

In 2017, Vaswani et al. introduced the Transformer architecture [29], bringing self-attention mechanisms to the forefront of deep learning. Since then, hybrid models such as AtteU-Net [30], Swin Transformer [31], and Tran-
sUNet [32] have emerged, combining CNN-based local feature extraction with the global contextual modeling capabilities of Transformers. Specifically, FAT-Net [33] effectively integrates the strengths of CNNs and Transformers for accurate melanoma segmentation. These hybrid approaches have established a new class of backbones for visual tasks, particularly in medical image segmentation [34].

2.2. SAM and MedSAM

Large Vision Models (LVMs) have played a critical role in advancing artificial intelligence in recent years [35]. As network architectures such as ResNet [36], U-Net [22], and ViT [37] continue to scale and evolve, large-scale pre-trained models have become increasingly influential in computer vision, particularly in image segmentation. Foundational models such as Mask2Former [38], Segmenter [39], and SETR [40], trained on large-scale datasets including ImageNet-21K [41] and JFT-300M [42], have demonstrated strong segmentation performance and generalization capabilities across diverse datasets.

Among these, the SAM has emerged as a representative foundation model for segmentation, notable for its strong zero-shot capability and cross-domain adaptability [5]. Its ability to adapt to varied image distributions and interpret complex visual patterns offers new opportunities for addressing challenging segmentation scenarios.

Despite substantial differences between natural and medical images, such foundation models have shown promising transferability to medical imaging tasks through domain-specific fine-tuning. Based on this idea, researchers developed MedSAM, a medical-specific foundation model built upon SAM and trained on over one million image–mask pairs from medical datasets. Quantitative results have validated its performance across various modalities and segmentation tasks, confirming its applicability in diverse clinical scenarios [6]. In addition, MedSAM has been shown to reduce annotation time by at least 82% compared to manual labeling. Among the prompt strategies originally designed for SAM, the bounding box prompt was found to provide clearer spatial context for regions of interest, thereby enabling more accurate identification of medical target areas. Given the strong generalization capabilities exhibited by SAM and MedSAM, across diverse imaging settings, the application of appropriate fine-tuning strategies holds considerable potential for enhancing their segmentation performance in domain-specific clinical tasks. As summarized in the Introduction, a growing body of research [7–16] has investigated various fine-tuning methodologies aimed at adapting these foundation models to specialized medical scenarios.

2.3. Parameter Fine-Tuning

Fine-tuning foundation models with domain-specific datasets enables task-specific adaptation and performance optimization. Full SFT involves

updating all model parameters, which often demands significant computational and hardware resources. To address this issue, a variety of parameter fine-tuning methods have been widely adopted. These techniques aim to reduce the number of trainable parameters while maintaining competitive performance, offering more efficient and flexible alternatives for adapting large-scale models [43].

Among these approaches, Prefix Tuning offers another lightweight strategy by prepending task-specific prefixes to the input sequence, avoiding modifications to the base model architecture [44]. In [45], Hu introduced LoRA, which injects low-rank updates into the attention matrices of Large Language Models (LLMs), significantly decreasing parameter counts. Adapter modules, first proposed by Houlsby et al. [46], introduce lightweight neural layers into models and have shown strong performance in domain adaptation scenarios. Both the Med-SA [9] and the LoRA-MedSAM [47] have since been validated on medical foundation models, where they provide efficient parameter tuning and competitive segmentation performance with only a small fraction of trainable weights.

Departing from original decoder substitution strategies typified by I-MedSAM [15], the parallel side network provides a non-intrusive, architecture-level enhancement pathway. By augmenting the backbone with a complementary functional module, it enhances task-specific representations while fully preserving the capabilities of the pretrained foundation model. This design ensures compatibility with pretrained weights and avoids alterations to the core architecture, thereby offering a flexible and non-intrusive enhancement mechanism [48].

3. Methods

This section begins with an overview of the MedSAM architecture, which forms the foundation of the proposed MedSAM-CA framework. As illustrated in Fig. 1, MedSAM-CA introduces several key enhancements to improve segmentation performance. A lightweight CNN branch, CBR-Net, runs in parallel with the ViT encoder to capture local spatial details. To effectively combine semantic and structural information, a multi-scale fusion module, Atte-FFB, integrates hierarchical features from both CBR-Net and the encoder. Furthermore, an Adapter module is incorporated into the encoder to strengthen global context modeling and enhance segmentation robustness.

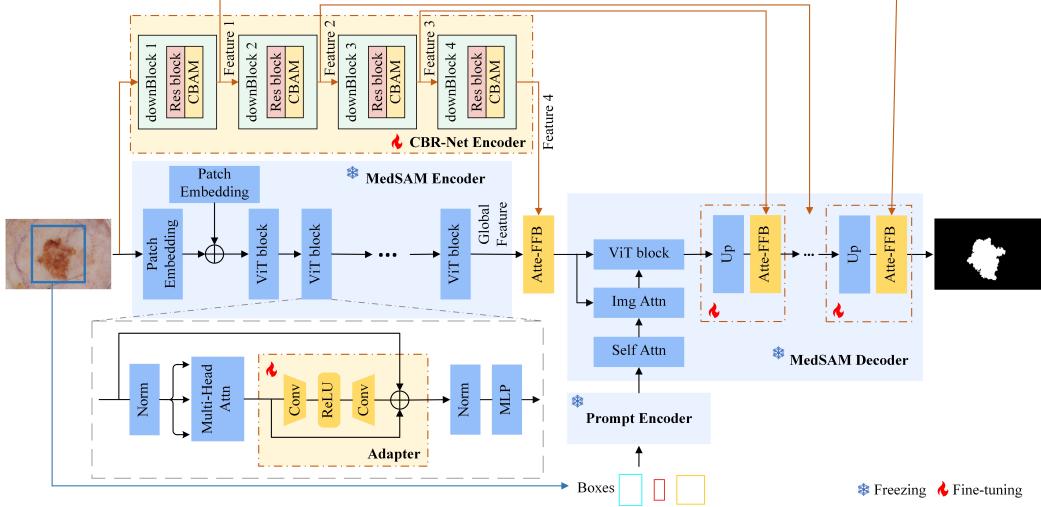


Figure 1: Overview of the proposed MedSAM-CA for medical image segmentation, including the parallel CBR-Net network, the multi-scale fusion module Atte-FFB, and the Adapter with MedSAM encoder.

3.1. Overview of MedSAM

MedSAM inherits the network architecture of SAM, as illustrated by the blue components in Fig. 1, and consists of three main modules: an image encoder, a prompt encoder, and a mask decoder. The image encoder, composed of 12 ViT blocks, is pre-trained using the Masked Autoencoder (MAE) strategy. Each ViT block contains a multi-head self-attention mechanism, a multi-layer perceptron (MLP), and normalization layers. The prompt encoder transforms user-defined prompts into contextual embeddings, which are subsequently fused with image features by the mask decoder through a cross-attention mechanism to generate precise segmentation masks. Notably, MedSAM effectively leverages bounding box prompts to guide attention toward target regions, thereby enhancing segmentation accuracy.

3.2. CBR-Net Encoder

This section presents the integration of CBR-Net into the MedSAM architecture. By connecting to the decoder via skip connections, CBR-Net enhances the delineation of fine boundaries through the incorporation of edge-aware features. Specifically, we design a four-layer CBR-Net that runs in parallel with the ViT encoder to extract fine-grained local representations. Using convolutional layers with progressively increasing receptive

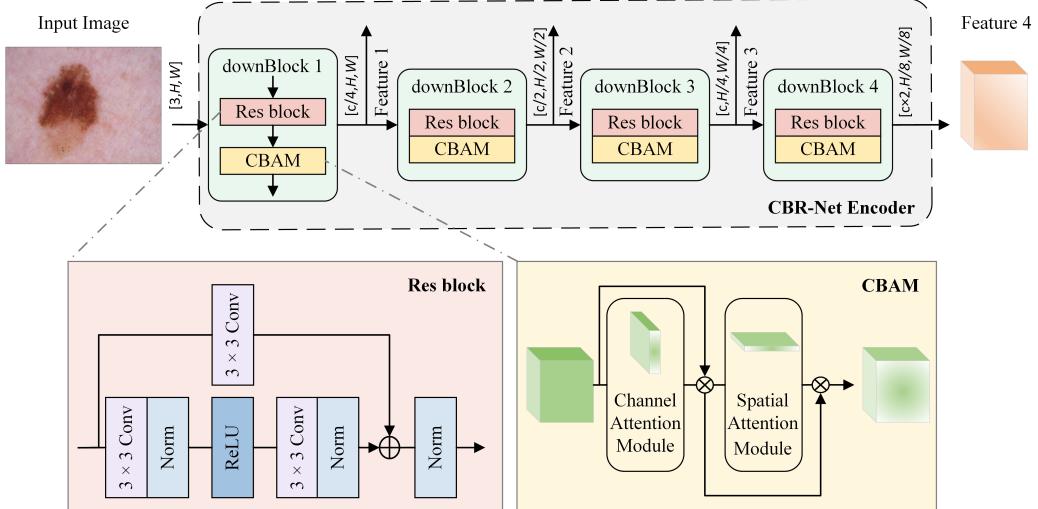


Figure 2: The parallel CBR-Net encoder consists of four cascaded downsampling blocks, each combining a residual block and a CBAM module.

fields, the network captures hierarchical and multi-scale spatial information. This architecture-level fine-tuning strategy equips the foundation model with enhanced representational capacity for boundary-sensitive segmentation.

The overall structure of CBR-Net, which consists of a cascade of residual blocks augmented by Convolutional Block Attention Modules (CBAM), is summarized in Fig. 2. Specifically, each residual block begins with a 3×3 convolution followed by a Rectified Linear Unit (ReLU) activation, and a second 3×3 convolution. The first 3×3 convolution simultaneously reduces spatial resolution by a factor of two and doubles the feature channel count, enabling the abstraction of higher-level semantic representations. To preserve identity information and mitigate feature degradation, a 1×1 convolution is applied in the shortcut path when discrepancies exist between input and output dimensions. The shortcut path is then combined with the main path output, forming the residual connection.

Each residual block is subsequently augmented with a CBAM, which applies channel-wise and spatial attention mechanisms to refine the feature representations. Detailed architectural descriptions of CBAM is available in [49].

Given an input image $x \in \mathbb{R}^{N \times H \times W}$, where N denotes the number of channels and $H \times W$ represents the spatial resolution, the input is pro-

cessed in parallel by the MedSAM ViT encoder and the CBR-Net. The branch fed into the ViT encoder produces a global feature map of dimensions $[c, H/16, W/16]$, where c denotes the number of feature channels and $H/16$ and $W/16$ indicate spatial dimensions following patch embedding and downsampling.

In the parallel branch, the four stages of CBR-Net progressively reduce the spatial resolution while increasing the feature dimensionality, generating hierarchical feature maps with resolutions of $[c/4, H, W]$, $[c/2, H/2, W/2]$, $[c, H/4, W/4]$, and $[c \times 2, H/8, W/8]$. We adopt a 4-stage configuration as it aligns naturally with common hierarchical CNN designs [22, 36], provides sufficient granularity for hierarchical representation, and ensures compatibility with the MedSAM ViT encoder’s downsampling ratio.

To comprehensively leverage multi-scale information, two complementary feature fusion strategies are introduced:

First, skip connections are established between the outputs of three CBR-Net stages (denoted as Features 1 to 3 in Fig. 1) and the corresponding layers in the MedSAM decoder. This hierarchical fusion strategy injects detailed local features at multiple spatial scales, thereby assisting the decoder in recovering fine-grained structures that may be lost in deeper Transformer layers. The implementation details of this fusion mechanism are described in Section 3.3. Second, the fourth-stage features from CBR-Net (Feature 4) and the Global Feature produced by the MedSAM encoder are aligned at matching spatial resolutions and then fused. This fusion enables direct interaction between local texture features and global semantic information, promoting more precise and context-aware segmentation results.

These strategies jointly enhance the decoder’s capacity to integrate global semantic information and fine-grained structural details, facilitating precise and context-aware segmentation.

3.3. Atte-FFB integration

In the original MedSAM architecture, the mask decoder produces the final segmentation output through a two-stage upsampling of high-dimensional feature maps. In contrast, during the decoding process, the proposed MedSAM-CA framework employs the Atte-FFB module to implement the two complementary feature fusion strategies introduced in the previous section. This design enables the integration of local details extracted by CBR-Net with the global contextual information encoded by the ViT, thereby improving mask quality and enhancing boundary precision.

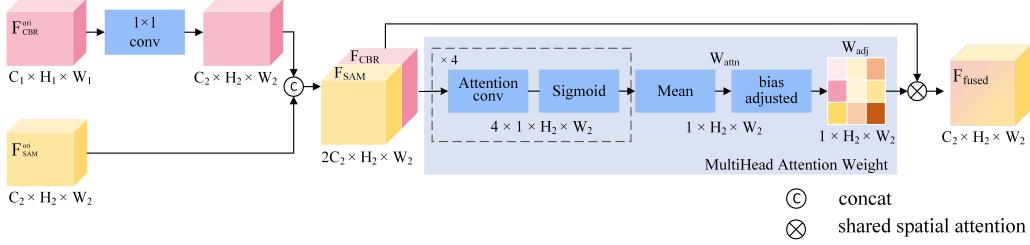


Figure 3: The network structure of the multi-scale fusion module Atte-FFB.

As illustrated in Fig. 3, we refer to the original features from CBR-Net as $F_{\text{CBR}}^{\text{ori}} \in \mathbb{R}^{C_1 \times H_1 \times W_1}$, and those from the MedSAM as $F_{\text{SAM}}^{\text{ori}} \in \mathbb{R}^{C_2 \times H_2 \times W_2}$. Feature fusion involves four sequential steps:

Feature Alignment. The original CBR-Net features $F_{\text{CBR}}^{\text{ori}}$ are first passed through a 1×1 convolution to align their channel dimensions with those of the corresponding MedSAM features $F_{\text{SAM}}^{\text{ori}}$. The resulting aligned features are denoted as F_{CBR} and F_{SAM} , which serve as inputs to the subsequent attention-based fusion.

Attention-Based Fusion Weight Generation. The aligned features F_{CBR} and F_{SAM} are concatenated to form a joint representation of size $[2C_2, H_2, W_2]$. This representation is passed through four parallel attention heads, each consisting of a convolutional layer followed by a Sigmoid activation. The resulting attention maps are averaged to produce a shared spatial attention weight $W_{\text{attn}} \in \mathbb{R}^{1 \times H_2 \times W_2}$, with values constrained to the range $[0, 1]$. This map reflects the relative importance of the two input streams at each spatial location.

Attention Weight Adjustment. To facilitate early-stage optimization of the randomly initialized CBR-Net, a learnable bias factor $b \in [0, 1]$ is introduced to initially emphasize its contribution in the fusion process. This mechanism, illustrated as the bias adjustment module in Fig. 3, is formulated as:

$$W_{\text{adj}} = W_{\text{attn}} \times (1 - b) + b. \quad (1)$$

This adjustment increases the weight of the CBR-Net features in the early training stages, allowing them to be learned more effectively. As training progresses, MedSAM features are gradually incorporated, enabling a balanced fusion of global and local information. This strategy also mitigates the adverse impact of early-stage noise from untrained CBR-Net parameters on the overall model convergence.

Weighted Feature Fusion. The final fused features F_{fused} are computed by element-wise combining the F_{SAM} and F_{CBR} using the adjusted attention weight:

$$F_{\text{fused}} = W_{\text{adj}} \cdot F_{\text{SAM}} + (1 - W_{\text{adj}}) \cdot F_{\text{CBR}}. \quad (2)$$

This fusion strategy incorporates global representations from the pretrained decoder along with local features extracted by CBR-Net, leading to improved training stability and enhanced segmentation accuracy.

3.4. MedSAM Encoder Adapter

To enhance the global feature extraction capability of MedSAM for two-dimensional (2D) medical image segmentation, we introduce an encoder-side adaptation strategy. Specifically, a standard Adapter layer [46] is embedded between the attention module and the MLP module within each ViT block of the image encoder, as illustrated in Fig. 1. This lightweight design enhances the model’s global feature extraction capabilities while avoiding the computational cost of full parameter fine-tuning. This approach has demonstrated effectiveness in improving segmentation accuracy for 2D medical images with reduced computational overhead.

4. Experiments

4.1. Datasets

To comprehensively evaluate the effectiveness and generalizability of the proposed MedSAM-CA framework, we conducted experiments on melanoma lesion segmentation and multi-organ abdominal segmentation. Melanoma segmentation is a challenging task in dermoscopy image analysis. The difficulty arises not only from lesion-specific characteristics—such as asymmetrical shapes, irregular boundaries, and heterogeneous pigmentation, but also from various confounding factors including hair occlusion, low contrast, skin tone variations, vascular structures, bubbles, artifacts, and suboptimal lighting conditions. These factors collectively complicate precise boundary delineation. Abdominal organ segmentation presents substantial challenges for boundary delineation and was therefore chosen as a representative evaluation task. Its difficulty arises from high inter-patient anatomical variability, blurred boundaries due to organ proximity, and intensity similarities between adjacent tissues such as the pancreas and fat in CT and Magnetic Resonance Imaging (MRI) scans.

Table 1: Dataset partition statistics for skin lesion and abdominal organ segmentation tasks.

	Skin Melanoma (images)			Abdominal Organs (cases)	
	ISIC 2016	ISIC 2017	ISIC 2018	AMOS 2022 CT	AMOS 2022 MRI
Train	900	2000	1815	200	40
Val	0	150	260	50	10
Test	379	600	519	50	10
Total	1279	2750	2594	300	60

Note: Numbers for ISIC datasets refer to 2D images, while those for AMOS 2022 dataset refer to 3D volumes.

Melanoma datasets. We conducted melanoma lesion segmentation experiments on three publicly available dermoscopic datasets: ISIC 2016 [50], ISIC 2017 [51], and ISIC 2018 [52]. These datasets include expert-annotated binary lesion masks and are widely used for benchmarking skin lesion segmentation models. The number of samples and data partitions for each dataset are summarized in Table 1. All input images were resized to 1024×1024 .

Abdominal organs datasets. For multi-organ segmentation, we used the AMOS 2022 dataset [53], which includes both CT and MRI scans with voxel-level annotations for 15 abdominal organs. The training set consists of 200 CT volumes and 40 MRI volumes, annotated with structures including the spleen, right kidney, left kidney, gallbladder, esophagus, liver, stomach, aorta, inferior vena cava, pancreas, right and left adrenal glands, duodenum, bladder, and prostate/uterus. The validation set comprises 100 CT and 20 MRI volumes, randomly split in a 1:1 ratio to form the validation and test subsets. Table 1 summarizes the dataset statistics and partitioning.

For preprocessing, intensity normalization is applied across all modalities. Specifically, CT volumes are clipped using a window width of 400 Hounsfield units (HU) and a window level of 40 HU. For MRI volumes, intensity values are clipped at the 0.5th and 99.5th percentiles of the intensity distribution. After clipping, all volumes are min–max normalized to the range $[0, 255]$. 3D CT and MRI volumes are processed on a per-slice basis, with each 2D slice resized to 1024×1024 pixels. Representative examples of input images and corresponding ground truth masks from the ISIC 2016-2018, AMOS 2022 CT, and AMOS 2022 MRI datasets are illustrated in Fig. 4.

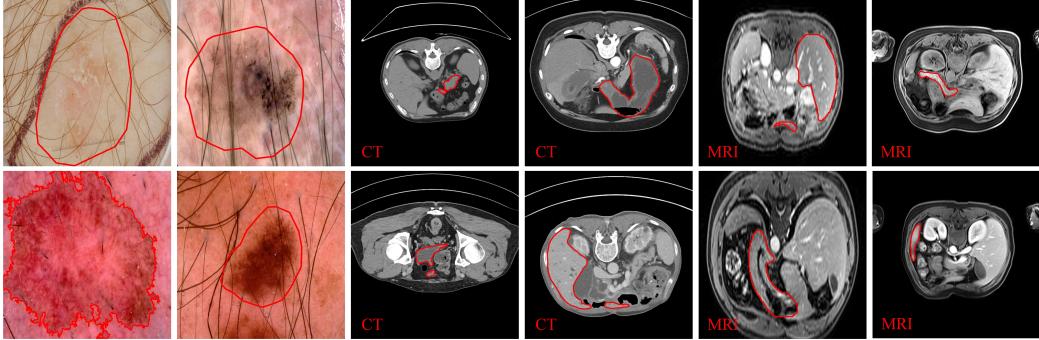


Figure 4: Representative examples from ISIC 2017 melanoma, AMOS 2022 CT, and AMOS 2022 MRI datasets.

4.2. Implementation details

The experiments were implemented on a Linux server equipped with an NVIDIA A100 GPU (80 GB VRAM) and Python 3.10.15. The deep learning framework used was PyTorch 2.5.1, running on Ubuntu 22.04.3 LTS with CUDA version 12.4. During training, the MedSAM-CA segmentation framework was optimized using the AdamW optimizer with a learning rate of 1×10^{-4} , a weight decay of 0.01, and a batch size of 8. The total number of training epochs was set to 80. The total loss function $\mathcal{L}_{\text{total}}$ was defined as a weighted sum of the binary cross-entropy loss \mathcal{L}_{BCE} and the Dice loss $\mathcal{L}_{\text{Dice}}$ as $\mathcal{L}_{\text{total}} = \alpha \cdot \mathcal{L}_{\text{BCE}} + (1 - \alpha) \cdot \mathcal{L}_{\text{Dice}}$, where the weighting factor α was set to 0.5. For melanoma segmentation, we used all available 4715 training images, comprising 900 from ISIC 2016, 2000 from ISIC 2017, and 1815 from ISIC 2018. The validation set included 410 images, drawn from 150 images in ISIC 2017 and 260 in ISIC 2018. For abdominal organ segmentation, the model was trained on 200 CT and 40 MRI volumes from the AMOS 2022 dataset, with validation performed on an additional 50 CT and 10 MRI volumes. Our code and models are available at: <https://github.com/Tianpeiting/MedSAM-CA.git>

4.3. Evaluation metrics

We evaluated segmentation performance using a combination of overlap-based and boundary-based metrics, including Dice coefficient (Dice), intersection over union (IoU), pixel accuracy (ACC), 95% Hausdorff Distance (HD95).

The Dice coefficient measures the overlap between the predicted segmentation A and the ground truth B , and is defined as:

$$\text{Dice} = \frac{2 \times |A \cap B|}{|A| + |B|}. \quad (3)$$

The IoU quantifies the ratio of the intersection to the union between the prediction and ground truth:

$$\text{IoU} = \frac{|A \cap B|}{|A \cup B|}. \quad (4)$$

The ACC calculates the ratio of correctly classified pixels (both foreground and background) to the total number of pixels:

$$\text{ACC} = \frac{TP + TN}{TP + TN + FP + FN}, \quad (5)$$

where TP , TN , FP , and FN denote true positives, true negatives, false positives, and false negatives, respectively.

All three metrics range from 0 to 1, with higher values indicating better segmentation quality.

The HD95 metric represents the 95th percentile of all minimum distances between boundary points of the predicted segmentation A and the ground truth B , and is defined as:

$$\text{HD95} = \max \{ P_{95} [d(a, S(B)) \mid a \in S(A)], P_{95} [d(b, S(A)) \mid b \in S(B)] \}, \quad (6)$$

where P_{95} denotes the 95th percentile operator, and $S(A)$ and $S(B)$ are the sets of boundary points for the predicted and ground truth masks, respectively. The distance functions are given by $d(a, S(B)) = \min_{b \in S(B)} \|a - b\|_2$ and $d(b, S(A)) = \min_{a \in S(A)} \|b - a\|_2$, where $S(A)$ and $S(B)$ denote the sets of boundary points for the predicted mask A and the ground truth mask B , respectively. HD95 is measured in pixels, with lower values indicating closer boundary alignment and better segmentation accuracy.

4.4. Comparison study on Melanoma Segmentation

In this section, we evaluate the performance of the proposed MedSAM-CA framework on melanoma segmentation. To validate our method, we conduct comparative analysis with traditional neural network-based segmentation approaches, such as UNet [22], nnU-Net [24], and DeepLabV3+ [27]. Several

MedSAM-based algorithms are also included in the comparison, including I-MedSAM [15], LoRA-MedSAM [47], and Med-SA [9]. Additionally, tailored algorithm for melanoma segmentation, such as FAT-Net [33] is also included.

For fairness, the bounding box was incorporated as an additional input channel during both training and inference for all baseline models that support multi-channel input, including UNet, nnU-Net, and DeepLabV3+. In contrast, MedSAM, I-MedSAM, LoRA-MedSAM, and Med-SA, which rely on prompt-based mechanisms, continued to receive the bounding box as a prompt input, consistent with their original implementation. In our study, bounding box prompts were generated from ground-truth masks by introducing random perturbations ranging from 0 to 20 pixels. The models were trained and evaluated using identical preprocessing pipelines, input resolutions, and dataset partitioning strategies to ensure consistency.

The melanoma segmentation performance of MedSAM-CA is evaluated on three publicly available skin lesion datasets: ISIC 2016, ISIC 2017, and ISIC 2018, and the quantitative results are detailed in Table 2. MedSAM-CA consistently outperforms all baseline methods across nearly all evaluation metrics on the three datasets. On ISIC 2016, MedSAM-CA achieves an ACC of 99.17%, an IoU of 97.18%, and a Dice score of 98.54%, representing absolute gains of 0.27%, 2.01%, and 0.62% over the second-best method, Med-SA. Furthermore, it yields reductions in boundary-related errors, with the HD95 of 10.96, compared to 25.56, achieved by the original MedSAM. These improvements underscore the enhanced boundary delineation capability of MedSAM-CA, particularly in cases involving subtle or irregular lesion margins.

Similar performance gains are observed on the ISIC 2017 dataset. MedSAM-CA achieves an HD95 of 18.56, reducing boundary deviations by 12.27, compared to MedSAM. It also attains an ACC of 98.71%, an IoU of 94.63%, and a Dice score of 97.10%, demonstrating consistent performance across varying imaging conditions.

It is also noteworthy that UNet consistently outperforms nnU-Net across the ISIC datasets. For example, on ISIC 2016, UNet achieves an ACC of 98.68% and a Dice score of 97.16%, slightly exceeding nnU-Net. One possible explanation is that UNet processes full-resolution inputs (1024×1024) during both training and inference, which may help retain global structural information and fine lesion boundaries more effectively than the patch-based reconstruction employed by nnU-Net.

On the ISIC 2018 dataset, MedSAM-CA maintains its superior perfor-

Table 2: Comparison of segmentation performance on ISIC 2016, ISIC 2017, and ISIC 2018 datasets.

ISIC 2016						
Model	Method	Journal/Year	HD95	ACC(%)	IoU(%)	Dice(%)
Specialized	FAT-Net [33]	MedIA 2022	17.93	97.79	91.54	95.45
Traditional	nnU-Net [24]	Nat. Methods 2021	17.35	98.44	94.20	96.98
	UNet [22]	MICCAI 2015	18.12	98.68	94.56	97.16
	DeepLabV3+ [27]	ECCV 2018	23.24	97.74	91.01	95.23
MedSAM based	MedSAM [6]	Nat. Commun. 2024	25.56	98.10	92.27	95.91
	I-MedSAM [15]	ECCV 2024	21.12	98.50	93.55	96.63
	LoRA-MedSAM [47]	MICAD 2024	16.77	98.62	94.80	97.30
	Med-SA [9]	MedIA 2025	<u>14.19</u>	98.90	<u>95.17</u>	<u>97.92</u>
	MedSAM-CA (Ours)	—	10.96	99.17	97.18	98.54
ISIC 2017						
Model	Method	Journal/Year	HD95	ACC(%)	IoU(%)	Dice(%)
Specialized	FAT-Net [33]	MedIA 2022	<u>19.60</u>	96.21	84.90	90.96
Traditional	nnU-Net [24]	Nat. Methods 2021	29.84	97.64	90.12	94.63
	UNet [22]	MICCAI 2015	30.05	98.06	90.86	95.05
	DeepLabV3+ [27]	ECCV 2018	33.73	96.93	86.66	92.65
MedSAM based	MedSAM [6]	Nat. Commun. 2024	30.83	97.73	88.73	93.83
	I-MedSAM [15]	ECCV 2024	30.10	97.78	90.13	94.63
	LoRA-MedSAM [47]	MICAD 2024	23.78	98.31	92.22	95.83
	Med-SA [9]	MedIA 2025	21.07	<u>98.60</u>	<u>93.40</u>	<u>96.49</u>
	MedSAM-CA (Ours)	—	18.56	98.71	94.63	97.10
ISIC 2018						
Model	Method	Journal/Year	HD95	ACC(%)	IoU(%)	Dice(%)
Specialized	FAT-Net [33]	MedIA 2022	<u>13.26</u>	98.31	91.92	95.24
Traditional	nnU-Net [24]	Nat. Methods 2021	22.18	98.52	92.78	96.20
	UNet [22]	MICCAI 2015	23.49	98.50	92.41	95.95
	DeepLabV3+ [27]	ECCV 2018	22.69	98.05	89.41	94.31
MedSAM based	MedSAM [6]	Nat. Commun. 2024	20.01	98.48	90.45	94.86
	I-MedSAM [15]	ECCV 2024	18.90	98.53	91.37	95.39
	LoRA-MedSAM [47]	MICAD 2024	14.18	98.74	94.07	97.02
	Med-SA [9]	MedIA 2025	14.85	<u>98.95</u>	<u>95.02</u>	<u>97.06</u>
	MedSAM-CA (Ours)	—	8.73	99.45	97.66	98.78

Note: The best-performing results are highlighted in **bold**, while the second-best results are underlined.

mance, achieving an HD95 of 8.73, an ACC of 99.45%, an IoU of 97.66%, and a Dice score of 98.78%. Although FAT-Net performs competitively on boundary-related metrics, achieving the second-best HD95 on both ISIC 2017 and ISIC 2018, its performance in region-based segmentation is relatively lim-

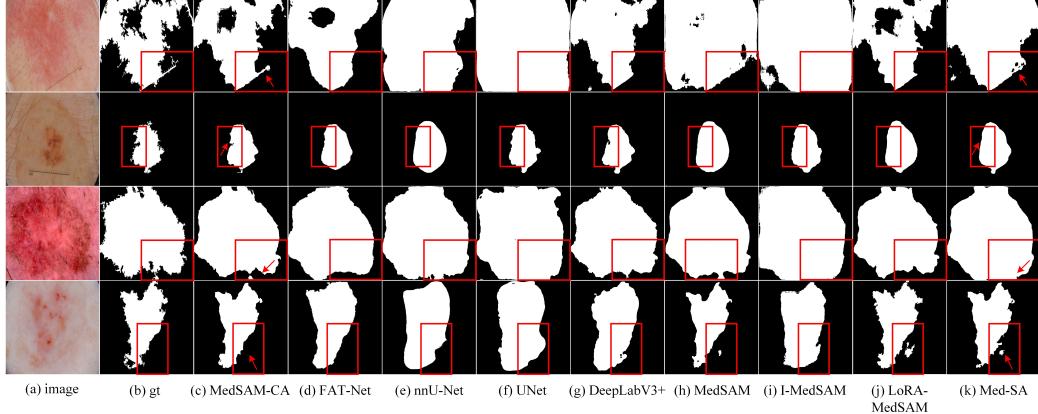


Figure 5: A visual comparison between MedSAM-CA and several competitive segmentation approaches on the ISIC 2017 dataset. (a) Input images. (b) Ground truth. (c) MedSAM-CA (Ours). (d) FAT-Net [33]. (e) nnU-Net [24]. (f) UNet [22]. (g) DeepLabV3+ [27]. (h) MedSAM [6]. (i) I-MedSAM [15]. (j) LoRA-MedSAM [47]. (k) Med-SA [9].

ited. For instance, on ISIC 2017, FAT-Net obtains an ACC of 96.21%, an IoU of 84.90%, and a Dice score of 90.96%, all of which fall below the scores achieved by the original MedSAM. These results highlight the importance of large-scale pretraining in improving segmentation accuracy and model generalizability.

Fig. 5 presents a visual comparison between MedSAM-CA and several representative segmentation methods on the ISIC 2017 dataset. As illustrated, MedSAM-CA consistently generates more accurate and complete segmentation masks, particularly in challenging cases with irregular boundaries (first row), varying lesion scales (second and third rows), and low contrast (fourth row). Compared to the second-best performer, Med-SA, MedSAM-CA exhibits superior boundary delineation and substantially fewer under-segmentation artifacts, as indicated by the arrows. These results underscore the robustness and effectiveness of MedSAM-CA in addressing difficult lesion segmentation scenarios.

Furthermore, Table 3 presents a comparison of model complexity and inference performance between MedSAM-CA and comparison methods, including total parameters, trainable parameters, inference time, and frames per second (FPS). It can be observed that our fine-tuning strategy remains lightweight overall. Compared to the original, untrained MedSAM, the average inference time of MedSAM-CA only increases by 0.0211 seconds, demon-

Table 3: Model complexity and inference efficiency comparison over ISIC 2016, ISIC 2017, and ISIC 2018 datasets.

Model	Method	Journal / Year	Params (M)	TrainParams (M)	Inference Time (s)	FPS
Specialized	FAT-Net [33]	MedIA 2022	30	30	0.0165	60.60
Traditional	nnU-Net [24]	Nat. Methods 2021	59	59	0.0258	38.75
	UNet [22]	MICCAI 2015	32	32	0.0136	73.52
	DeepLabV3+ [27]	ECCV 2018	24	24	0.0128	78.12
MedSAM based	MedSAM [6]	Nat. Commun. 2024	94	0	0.0959	10.43
	I-MedSAM [15]	ECCV 2024	96	2	0.2240	4.46
	LoRA-MedSAM [47]	MICAD 2024	96	2	0.1035	9.66
	Med-SA [9]	MedIA 2025	101	7	0.1206	8.29
	MedSAM-CA (Ours)	—	101	7	0.1170	8.55

Note: Inference time is computed as the average over 600 forward passes, measured after sufficient model warm-up. The input resolution is kept consistent with training (1024×1024), and the batch size is set to 1.

strating high efficiency despite the added fine-tuning modules.

These results further validate the reliability of combining foundational models with lightweight structural adaptations for high-precision medical image segmentation, especially in tasks requiring fine boundary delineation and strong generalization.

4.5. Comparison study on Abdominal Organs Segmentation

To evaluate the generalization ability of MedSAM-CA across imaging modalities, we conduct experiments on the AMOS 2022 dataset, which comprises both CT and MRI scans. The evaluation follows the same comparison protocol as described in Section 4.4, which includes traditional neural network-based segmentation models, MedSAM-based models, and Hermes-R [28], a method specifically designed for abdominal organ segmentation. The quantitative results are summarized in Table 4.

For CT modality, MedSAM-CA achieves the highest segmentation accuracy across all metrics. Specifically, it reaches a Dice score of 93.32% and an IoU of 88.13%, outperforming the nnU-Net by 3.79% and 3.29%, respectively. In terms of boundary precision, it achieves an HD95 of 7.22, indicating more precise localization of organ contours compared to all other methods. These results confirm the model’s effectiveness in capturing structural information from high-resolution CT data.

For MRI modality, the improvements are even more pronounced. MedSAM-CA achieves a Dice score of 94.51%, an IoU of 90.10%, and an HD95 of 7.75. These performance gains can be partially attributed to the

Table 4: Comparison of segmentation performance on AMOS 2022 dataset.

Model	Method	Journal / Year	Modality	HD95	ACC(%)	IoU(%)	Dice(%)
Specialized	Hermes-R [28]	CVPR 2024	CT	13.67	99.73	84.09	89.14
			MRI	16.34	99.69	76.67	86.24
Traditional	nnU-Net [24]	Nat. Methods 2021	CT	8.65	99.78	84.84	89.53
			MRI	9.16	99.70	84.17	88.89
	UNet [22]	MICCAI 2015	CT	<u>7.64</u>	99.80	85.37	90.72
	DeepLabV3+ [27]	ECCV 2018	MRI	<u>8.80</u>	99.71	84.76	90.84
			CT	16.18	99.78	81.31	88.76
			MRI	23.66	99.55	81.39	88.94
MedSAM based	MedSAM [6]	Nat. Commun. 2024	CT	14.81	99.70	75.04	84.12
			MRI	12.22	99.59	80.80	88.45
	I-MedSAM [15]	ECCV 2024	CT	13.27	99.22	84.51	91.19
	LoRA-MedSAM [47]	MICAD 2024	MRI	12.48	98.30	84.66	91.40
			CT	<u>7.82</u>	<u>99.89</u>	<u>87.36</u>	92.81
	Med-SA [9]	MedIA 2025	MRI	<u>8.32</u>	<u>99.84</u>	89.51	<u>94.12</u>
	MedSAM-CA (Ours)	—	CT	7.76	99.86	87.01	<u>93.26</u>
			MRI	<u>8.12</u>	99.82	<u>89.69</u>	94.05
	MedSAM-CA (Ours)	—	CT	7.22	99.90	88.13	93.32
			MRI	7.75	99.85	90.10	94.51

Note: The best-performing results are highlighted in **bold**, while the second-best results are underlined.

training distribution of MedSAM, in which MRI accounts for the largest proportion. The prevalence of MRI in pretraining enhances the model’s prior knowledge for this imaging modality, contributing to its superior performance in handling lower overall image contrast and the increased complexity of soft-tissue structures typically encountered in MRI scans.

Hermes-R, developed specifically for abdominal organ segmentation, demonstrates competitive performance in region-based metrics on CT scans, achieving a Dice score of 89.14% and an IoU of 84.09%. However, its boundary accuracy remains limited, with an HD95 of 13.67, which is less favorable compared to MedSAM-CA. This discrepancy may stem from the lack of large-scale pretraining in Hermes-R, which potentially limits its ability to generalize across organ shapes and image variations. In contrast, MedSAM-CA builds on the advantages of pretrained foundational models and incorporates targeted structural adaptations, resulting in more accurate segmentation at both the region and boundary levels. These comparisons underscore the practical value of combining generalizable pretrained representations with lightweight architecture-level fine-tuning for medical image segmentation tasks involving diverse modalities.

Fig. 6 presents a visual comparison of segmentation results on the AMOS

2022 dataset across multiple representative methods. The left two columns present CT examples, while the right two columns show MRI cases.

As illustrated, MedSAM-CA consistently produces more accurate and coherent segmentation masks across both modalities. In particular, the second and third columns highlight the left kidney and the duodenum, which are both anatomically challenging structures to segment. The left kidney is difficult to delineate due to its proximity to tissues with similar intensities, such as the spleen and muscle, as well as variations in its shape and position across patients. The duodenum presents additional challenges because of its thin wall, curved tubular structure, and high anatomical variability, which often lead to confusion with adjacent organs including the pancreas and intestinal loops. By effectively integrating local and global features, MedSAM-CA enables more precise boundary delineation for both large anatomical structures and small, low-contrast regions. Notably, the organ highlighted by the blue mask in the first and fourth columns is the stomach. Accurate boundary delineation for these organs remains challenging due to their variable shapes, intensity similarities with surrounding tissues, and indistinct edges in CT/MRI images. Compared with the original MedSAM, which often results in incomplete or blurred boundaries, especially for large or poorly contrasted organs, MedSAM-CA provides more consistent contours. These results show the advantage of incorporating both global and local feature enhancements in the proposed framework.

In addition, UNet generally preserves overall organ structure, it tends to miss details in smaller or less distinct regions (first column). DeepLabV3+, on the other hand, struggles with segmenting several organs accurately (see the first and third columns), especially in the MRI modality. This limitation likely stems from its original design for natural image segmentation, where texture continuity and edge sharpness differ markedly from the low-contrast and noise-prone characteristics of medical images. These observations are consistent with the quantitative results and further demonstrate the robustness and adaptability of the proposed MedSAM-CA across diverse imaging conditions.

In summary, MedSAM-CA integrates the representational advantages of foundation model pretraining with task-specific structural adaptations in a unified framework. Its consistent improvements across both region-based and boundary-aware metrics demonstrate the value of jointly leveraging global semantic and local spatial features, especially in challenging medical image segmentation tasks.

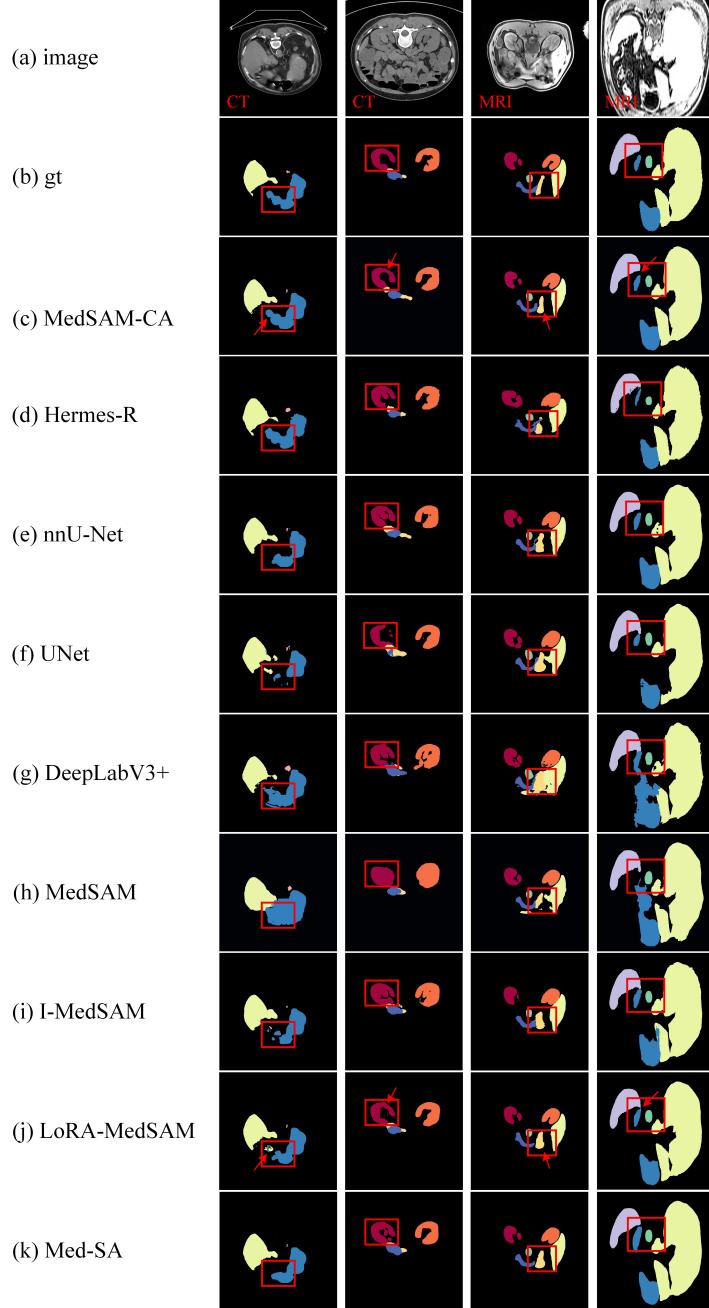


Figure 6: A visual comparison between MedSAM-CA and several competitive segmentation approaches on the AMOS 2022 dataset. (a) Input images. (b) Ground truth. (c) MedSAM-CA (Ours). (d) Hermes-R [28]. (e) nnU-Net [24]. (f) UNet [22]. (g) DeepLabV3+ [27]. (h) MedSAM [6]. (i) I-MedSAM [15]. (j) LoRA-MedSAM [47]. (k) Med-SA [9].

Table 5: Model Components Ablation Study on ISIC 2018 and AMOS 2022

Modality	Method	HD95	ACC (%)	IoU (%)	Dice (%)	TrainParams (M)
Dermoscopy	Full (CBR-Net+Atte-FFB+Adapter)	8.73	99.45	97.66	98.78	6.8
	w/o Adapter	11.00	99.29	96.63	98.26	5.4
	w/o Atte-FFB	14.83	99.09	95.14	97.48	5.2
	w/o CBR-Net & Atte-FFB	16.52	98.92	94.37	97.05	1.4
	MedSAM (baseline)	20.01	98.48	90.45	94.86	0
CT	Full (CBR-Net+Atte-FFB+Adapter)	7.22	99.90	88.13	93.32	6.8
	w/o Adapter	8.14	99.89	87.53	92.94	5.4
	w/o Atte-FFB	8.37	99.86	86.96	92.45	5.2
	w/o CBR-Net & Atte-FFB	8.45	99.85	86.52	92.18	1.4
	MedSAM (baseline)	14.82	99.70	75.04	84.12	0
MRI	Full (CBR-Net+Atte-FFB+Adapter)	7.75	99.85	90.10	94.51	6.8
	w/o Adapter	7.93	99.85	90.01	94.45	5.4
	w/o Atte-FFB	8.08	99.82	89.84	94.01	5.2
	w/o CBR-Net & Atte-FFB	8.15	99.80	89.35	93.69	1.4
	MedSAM (baseline)	12.22	99.59	80.80	88.45	0

Note: “w/o” indicates the removal of specific components. “w/o Adapter” disables the Adapter module; “w/o Atte-FFB” removes the multi-scale fusion block; “w/o CBR-Net & Atte-FFB” removes both the CBR-Net branch and fusion. “MedSAM (baseline)” refers to the original, unfine-tuned MedSAM model.

4.6. Ablation study

To validate the contribution of each component within the MedSAM-CA framework, we conducted an ablation study across three imaging modalities: dermoscopy (ISIC 2018), CT (AMOS 2022), and MRI (AMOS 2022). The ablation experiments focus on three critical modules: the parallel CNN branch (CBR-Net), the feature fusion block (Atte-FFB), and the Adapter module.

Since the Atte-FFB module operates in conjunction with CBR-Net, we cannot keep the Atte-FFB module only and remove CBR-Net module. We conduct the ablation experiment of “w/o CBR-Net & Atte-FFB”. After CBR-Net module, we substitute the Atte-FFB module with simple element-wise addition to maintain the architectural consistency. Therefore, we conduct the ablation experiment of “w/o Atte-FFB”.

Table 5 summarizes the ablation results. The results confirm the effectiveness of each proposed component. On dermoscopy images, removing the CBR-Net and Atte-FFB module from the full MedSAM-CA framework leads to a noticeable performance drop, with HD95 increasing from 8.73 to 16.52 and IoU decreasing from 97.66% to 94.37%. This decline suggests that the

parallel CNN branch, CBR-Net, provides complementary local features to the ViT encoder, potentially contributing to improved boundary segmentation accuracy.

After removing the Adapter module, the HD95 and IoU degrade to 11.00 and 96.63%, respectively, suggesting that lightweight tuning of global features provides additional benefits for segmentation. In comparison, ablating only the Atte-FFB module leads to a larger performance drop, with HD95 and IoU decreasing to 14.83 and 95.14%. These results indicate that Atte-FFB plays a more important role in the proposed framework, as effective multi-scale feature fusion enhances the model’s ability to integrate local and global representations.

Experimental results on CT and MRI align with those on dermoscopy, as presented in Table 5. The TrainParams column indicates that fine-tuning a small subset of parameters can effectively enhance segmentation performance, underscoring the importance of lightweight architecture-level fine-tuning in improving accuracy.

4.7. Experiments on Limited Training Dataset

To evaluate the effectiveness of MedSAM-CA framework in reducing reliance on training data, we conducted experiments on limited data settings. In all experiments, the validation and test sets remained consistent with those defined in the original challenge splits, as shown in Table 1, ensuring fair comparisons across different training regimes.

Table 6 presents a comparison of MedSAM-CA, UNet, and I-MedSAM on melanoma segmentation using different training sample sizes, selected to reflect performance in data-scarce settings. With only 2% of the original training data, UNet achieved a Dice of 90.27%. In contrast, MedSAM-CA achieved a Dice of 94.43%, reaching 97.25% (94.43%/97.10%) of its full-data performance, demonstrating the strong generalization capability of large foundation models. Notably, I-MedSAM achieved a Dice of 91.39%, which did not surpass the original MedSAM (93.83% as shown in Table 2). This may be attributed to the complete replacement of MedSAM’s original decoder in I-MedSAM, which may require additional training data to be effectively retrained for strong segmentation results. In addition, MedSAM-CA achieves a Dice score of 93.91% and an IoU of 89.10% using only 0.4% of annotated training data. By contrast, UNet achieves a Dice score of 87.10% and an IoU of 78.47% under the same setting.

Table 6: Segmentation performance of UNet, I-MedSAM and MedSAM-CA on limited training data on ISIC 2017.

Dataset	Method	Evaluation (%)	Training Samples			
			0.4%	1.2%	2%	100%
ISIC 2017	UNet	IoU	78.47	82.26	83.04	90.86
		Dice	87.10	89.60	90.27	95.05
	I-MedSAM	IoU	81.01	83.58	84.59	90.13
		Dice	89.20	90.79	91.39	94.63
	MedSAM-CA	IoU	89.10	89.61	89.75	94.63
		Dice	93.91	94.34	94.43	97.10

Table 7: Segmentation performance of UNet, I-MedSAM and MedSAM-CA on limited training data on AMOS 2022.

Dataset	Method	Evaluation (%)	Training Samples			
			2%	6%	10%	100%
AMOS 2022 CT	UNet	IoU	81.21	83.46	83.99	85.37
		Dice	88.54	90.12	90.51	90.72
	I-MedSAM	IoU	70.75	73.84	75.55	84.51
		Dice	81.19	83.61	84.60	91.19
	MedSAM-CA	IoU	85.33	85.79	86.86	88.13
		Dice	91.47	91.78	92.50	93.32
AMOS 2022 MRI	UNet	IoU	82.01	83.23	83.97	84.76
		Dice	89.26	90.13	90.64	90.84
	I-MedSAM	IoU	75.40	78.83	80.03	84.66
		Dice	84.91	87.40	88.21	91.40
	MedSAM-CA	IoU	88.52	89.06	89.57	90.10
		Dice	93.52	93.83	94.22	94.51

For the AMOS 2022 dataset, the segmentation performance of MedSAM-CA, UNet, and I-MedSAM across different training sample sizes is summarized in Table 7. Some of the training sample settings were inspired by recent multi-organ segmentation studies [16]. MedSAM-CA achieved an IoU of 89.57% on MRI using only 10% of the full annotated CT and MRI data, reaching 99.41% of the full-data performance (89.57%/90.10%). In contrast, I-MedSAM achieved only 80.03% IoU under the same conditions, similar to the MedSAM (80.80% as shown in Table 4). These results demonstrate the strong efficiency and clinical applicability of MedSAM-CA in data-limited

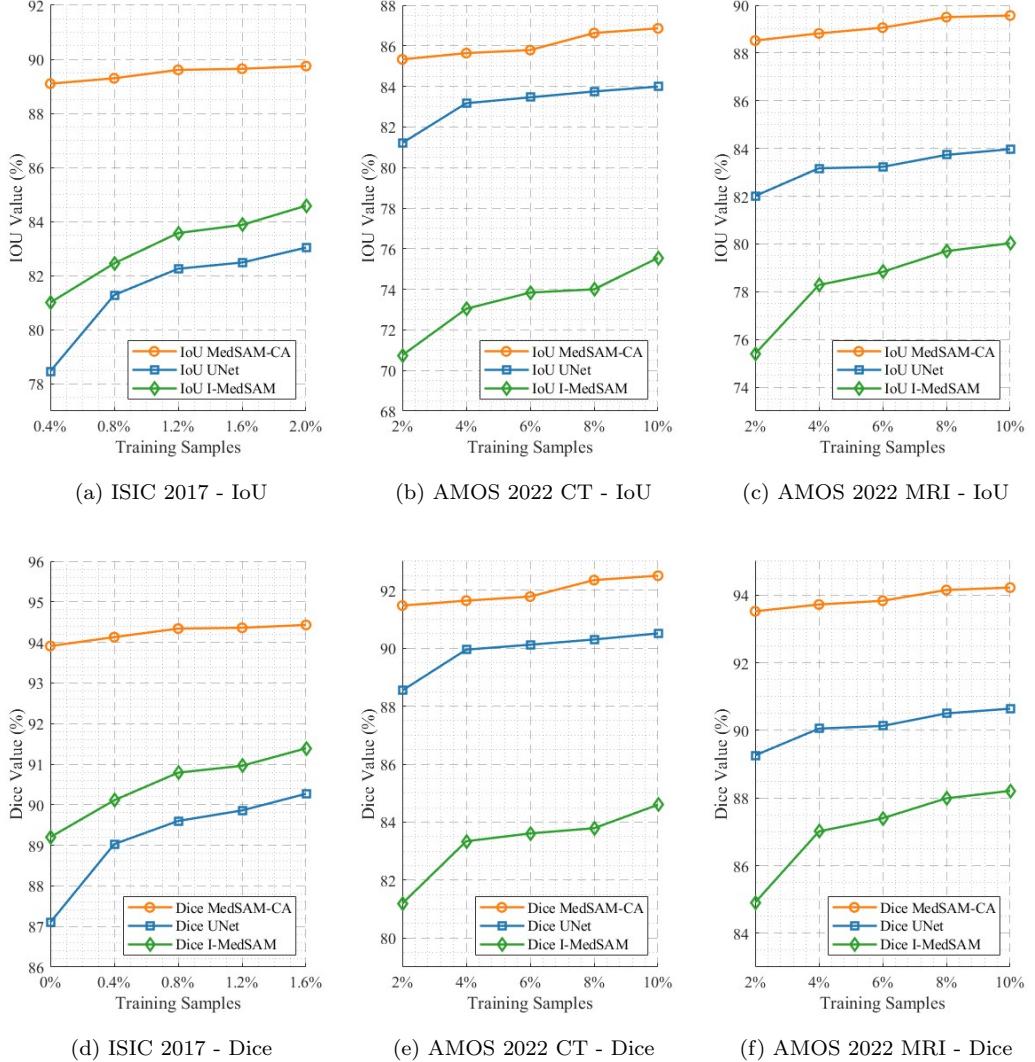


Figure 7: Performance comparison between MedSAM-CA, UNet and I-MedSAM on varying training sample sizes across three datasets. (a)-(c) IoU performance on ISIC 2017, AMOS 2022 CT, and AMOS 2022 MRI. (d)-(f) Dice performance on the same datasets.

scenarios.

Fig. 7 presents how all models improve in segmentation accuracy with more training samples. However, MedSAM-CA consistently surpasses UNet and I-MedSAM across all data scales, demonstrating strong robustness in low-resource scenarios. Moreover, for both IoU and Dice evaluation scores,

MedSAM-CA shows a greater performance gain over UNet in the MRI modality compared to CT. This advantage may be attributed to the higher proportion of MRI data in the original MedSAM’s pretraining distribution. These findings validate the effectiveness of adapting foundation models to limited-data medical segmentation tasks and highlight the strong segmentation capability of MedSAM-CA.

5. Discussion and conclusion

In this work, we proposed MedSAM-CA, a lightweight architecture-level fine-tuning of the MedSAM framework designed to enhance segmentation performance for challenging medical imaging tasks. MedSAM-CA incorporates two complementary architectural innovations, CBR-Net and Atte-FFB. CBR-Net is a lightweight CNN branch operating in parallel with the ViT encoder to enhance the capture of local textures and boundary-specific details. The extracted features from CBR-Net are connected to the decoder via skip connections, effectively preserving spatial details critical for accurate boundary delineation. Atte-FFB is a feature fusion structure that leverages attention mechanisms to robustly integrates hierarchical features from both the convolutional and Transformer-based branches, thereby enabling more effective multi-scale representation alignment. Finally, an Adapter module was embedded within the encoder to facilitate efficient fine-tuning of global representations without modifying the original pretrained parameters.

Our experimental results confirm that MedSAM-CA improves segmentation performance in scenarios characterized by subtle and irregular boundaries, such as melanoma and abdominal organ segmentation. Importantly, MedSAM-CA achieves these improvements with only 0.0211 seconds increase in inference time per image. Furthermore, ablation studies confirmed the individual contributions of each proposed module. The results from limited-data experiments also demonstrate that MedSAM-CA effectively reduces dependence on large-scale annotated datasets, highlighting its potential for application in clinical scenarios where data accessibility is restricted due to privacy concerns or annotation costs.

Despite these advances, several limitations warrant consideration. First, our investigation was restricted to the publicly available ViT-B-based MedSAM model, as larger-scale ViT architectures (e.g., ViT-L or ViT-H) are currently unavailable within the MedSAM repository. Second, our experiments primarily utilized bounding-box prompts with perturbations, following prior

recommendations for MedSAM. In future work, more robust prompt types such as point-based inputs, user-drawn strokes, or other interactive modalities may be explored.

In conclusion, the proposed MedSAM-CA framework represents an effective and computationally efficient approach to adapting foundational segmentation models for precise boundary delineation in medical images. The proposed framework offers a practical solution for precise medical segmentation, improves diagnostic consistency, and provides a foundation for developing advanced computer-aided diagnosis systems. Future studies may consider further integration of MedSAM-CA with advanced architectural components, exploration of diverse prompting strategies, and validation on additional clinical imaging tasks to further extend its applicability.

References

- [1] S. Bakas, H. Akbari, A. Sotiras, M. Bilello, M. Rozycki, J. S. Kirby, J. B. Freymann, K. Farahani, C. Davatzikos, Advancing the cancer genome atlas glioma mri collections with expert segmentation labels and radiomic features, *Scientific data* 4 (1) (2017) 1–13.
- [2] K. Doi, Computer-aided diagnosis in medical imaging: historical review, current status and future potential, *Computerized medical imaging and graphics* 31 (4-5) (2007) 198–211.
- [3] S. Pouyanfar, S. Sadiq, Y. Yan, H. Tian, Y. Tao, M. P. Reyes, M.-L. Shyu, S.-C. Chen, S. S. Iyengar, A survey on deep learning: Algorithms, techniques, and applications, *ACM computing surveys (CSUR)* 51 (5) (2018) 1–36.
- [4] H. H. Lee, Y. Gu, T. Zhao, Y. Xu, J. Yang, N. Usuyama, C. Wong, M. Wei, B. A. Landman, Y. Huo, et al., Foundation models for biomedical image segmentation: A survey, *arXiv preprint arXiv:2401.07654* (2024).
- [5] A. Kirillov, E. Mintun, N. Ravi, H. Mao, C. Rolland, L. Gustafson, T. Xiao, S. Whitehead, A. C. Berg, W.-Y. Lo, et al., Segment anything, in: *Proceedings of the IEEE/CVF international conference on computer vision*, 2023, pp. 4015–4026. doi:10.1109/ICCV51070.2023.00371.

- [6] J. Ma, Y. He, F. Li, L. Han, C. You, B. Wang, Segment anything in medical images, *Nature Communications* 15 (1) (2024) 654. doi:<https://doi.org/10.1038/s41467-024-44824-z>.
- [7] H. Zhicheng, W. Yipeng, L. Xiao, Deep learning-based detection of impacted teeth on panoramic radiographs, *Biomedical Engineering and Computational Biology* 15 (2024) 11795972241288319.
- [8] K. Pang, S. Zeng, Mediastinal image detection and segmentation based on medSAM, in: Fourth International Conference on Image Processing and Intelligent Control (IPIC 2024), Vol. 13250, SPIE, 2024, pp. 145–155.
- [9] J. Wu, Z. Wang, M. Hong, W. Ji, H. Fu, Y. Xu, M. Xu, Y. Jin, Medical SAM adapter: Adapting segment anything model for medical image segmentation, *Medical image analysis* 102 (2025) 103547.
- [10] K. Li, P. Rajpurkar, Adapting segment anything models to medical imaging via fine-tuning without domain pretraining, in: AAAI 2024 Spring Symposium on Clinical Foundation Models, 2024.
- [11] J. Wu, M. Xu, One-prompt to segment all medical images, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2024, pp. 11302–11312.
- [12] Q. Wu, Y. Zhang, M. Elbatel, Self-prompting large vision models for few-shot medical image segmentation, in: MICCAI workshop on domain adaptation and representation transfer, Springer, 2023, pp. 156–167.
- [13] Y. Tang, D. Yang, W. Li, H. R. Roth, B. Landman, D. Xu, V. Nath, A. Hatamizadeh, Self-supervised pre-training of Swin Transformers for 3D medical image analysis, in: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, 2022, pp. 20730–20740.
- [14] W. Lei, W. Xu, K. Li, X. Zhang, S. Zhang, MedLSAM: Localize and segment anything model for 3D CT images, *Medical Image Analysis* 99 (2025) 103370.
- [15] X. Wei, J. Cao, Y. Jin, M. Lu, G. Wang, S. Zhang, I-medSAM: Implicit medical image segmentation with segment anything, in: European Conference on Computer Vision, Springer, 2024, pp. 90–107.

- [16] Y. Li, B. Jing, Z. Li, J. Wang, Y. Zhang, Plug-and-play segment anything model improves nnUNet performance, *Medical physics* (2025).
- [17] Y.-T. Hsiao, C.-L. Chuang, J.-A. Jiang, C.-C. Chien, A contour based image segmentation algorithm using morphological edge detection, in: 2005 IEEE International Conference on systems, man and cybernetics, Vol. 3, 2005, pp. 2962–2967.
- [18] T. Heimann, H.-P. Meinzer, Statistical shape models for 3D medical image segmentation: a review, *Medical image analysis* 13 (4) (2009) 543–563.
- [19] H. G. Kaganami, B. Zou, Region-based segmentation versus edge detection, in: 2009 Fifth International Conference on Intelligent Information Hiding and Multimedia Signal Processing, IEEE, 2009, pp. 1217–1221.
- [20] E. E. Nithila, S. Kumar, Segmentation of lung from CT using various active contour models, *Biomedical Signal Processing and Control* 47 (2019) 57–62.
- [21] J. Long, E. Shelhamer, T. Darrell, Fully convolutional networks for semantic segmentation, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2015, pp. 3431–3440.
- [22] O. Ronneberger, P. Fischer, T. Brox, U-net: Convolutional networks for biomedical image segmentation, in: *Medical image computing and computer-assisted intervention–MICCAI 2015: 18th international conference, Munich, Germany, October 5–9, 2015, proceedings, part III* 18, Springer, 2015, pp. 234–241.
- [23] N. P. Bindu, P. N. Sastry, RETRACTED ARTICLE: Automated brain tumor detection and segmentation using modified UNet and ResNet model, *Soft Computing* 27 (13) (2023) 9179–9189.
- [24] F. Isensee, P. F. Jaeger, S. A. Kohl, J. Petersen, K. H. Maier-Hein, nnU-Net: A self-configuring method for deep learning-based biomedical image segmentation, *Nature methods* 18 (2) (2021) 203–211. doi: <https://doi.org/10.1038/s41592-020-01008-z>.

- [25] F. Milletari, N. Navab, S.-A. Ahmadi, V-net: Fully convolutional neural networks for volumetric medical image segmentation, in: 2016 fourth international conference on 3D vision (3DV), IEEE, 2016, pp. 565–571.
- [26] Ö. Çiçek, A. Abdulkadir, S. S. Lienkamp, T. Brox, O. Ronneberger, 3D U-Net: learning dense volumetric segmentation from sparse annotation, in: Medical Image Computing and Computer-Assisted Intervention—MICCAI 2016: 19th International Conference, Athens, Greece, October 17-21, 2016, Proceedings, Part II 19, Springer, 2016, pp. 424–432.
- [27] L.-C. Chen, Y. Zhu, G. Papandreou, F. Schroff, H. Adam, Encoder-decoder with atrous separable convolution for semantic image segmentation, in: Proceedings of the European conference on computer vision (ECCV), Vol. 11211, Munich, Germany, 2018, pp. 801–818. doi: [10.1007/978-3-030-01234-2_49](https://doi.org/10.1007/978-3-030-01234-2_49).
- [28] Y. Gao, Training like a medical resident: Context-prior learning toward universal medical image segmentation, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2024, pp. 11194–11204.
- [29] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, I. Polosukhin, Attention is all you need, Advances in neural information processing systems 30 (2017).
- [30] J. Schlemper, O. Oktay, M. Schaap, M. Heinrich, B. Kainz, B. Glocker, D. Rueckert, Attention gated networks: Learning to leverage salient regions in medical images, Medical image analysis 53 (2019) 197–207. doi:<https://doi.org/10.1016/j.media.2019.01.012>.
- [31] Z. Liu, Y. Lin, Y. Cao, H. Hu, Y. Wei, Z. Zhang, S. Lin, B. Guo, Swin transformer: Hierarchical vision transformer using shifted windows, in: Proceedings of the IEEE/CVF international conference on computer vision, 2021, pp. 10012–10022.
- [32] J. Chen, Y. Lu, Q. Yu, X. Luo, E. Adeli, Y. Wang, L. Lu, A. L. Yuille, Y. Zhou, Transunet: Transformers make strong encoders for medical image segmentation, arXiv preprint arXiv:2102.04306 (2021).
- [33] H. Wu, S. Chen, G. Chen, W. Wang, B. Lei, Z. Wen, FAT-Net: Feature adaptive transformers for automated skin lesion segmentation, Medical

image analysis 76 (2022) 102327. doi:<https://doi.org/10.1016/j.media.2021.102327>.

- [34] K. Lu, Y. Xu, Y. Yang, Comparison of the potential between transformer and cnn in image classification, in: ICMLCA 2021; 2nd International Conference on Machine Learning and Computer Application, VDE, 2021, pp. 1–6.
- [35] J. Wang, Z. Liu, L. Zhao, Z. Wu, C. Ma, S. Yu, H. Dai, Q. Yang, Y. Liu, S. Zhang, et al., Review of large vision models and visual prompt engineering, *Meta-Radiology* 1 (3) (2023) 100047.
- [36] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image recognition, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2016, pp. 770–778.
- [37] L. Yuan, Y. Chen, T. Wang, W. Yu, Y. Shi, Z.-H. Jiang, F. E. Tay, J. Feng, S. Yan, Tokens-to-token vit: Training vision transformers from scratch on imagenet, in: Proceedings of the IEEE/CVF international conference on computer vision, 2021, pp. 558–567.
- [38] B. Cheng, I. Misra, A. G. Schwing, A. Kirillov, R. Girdhar, Masked-attention mask transformer for universal image segmentation, in: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, 2022, pp. 1290–1299.
- [39] R. Strudel, R. Garcia, I. Laptev, C. Schmid, Segmenter: Transformer for semantic segmentation, in: Proceedings of the IEEE/CVF international conference on computer vision, 2021, pp. 7262–7272.
- [40] S. Zheng, J. Lu, H. Zhao, X. Zhu, Z. Luo, Y. Wang, Y. Fu, J. Feng, T. Xiang, P. H. Torr, et al., Rethinking semantic segmentation from a sequence-to-sequence perspective with transformers, in: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, 2021, pp. 6881–6890.
- [41] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, F.-F. Li, Imagenet: A large-scale hierarchical image database, in: 2009 IEEE conference on computer vision and pattern recognition, IEEE, 2009, pp. 248–255.

- [42] C. Sun, A. Shrivastava, S. Singh, A. Gupta, Revisiting unreasonable effectiveness of data in deep learning era, in: Proceedings of the IEEE international conference on computer vision, 2017, pp. 843–852.
- [43] N. Ding, Y. Qin, G. Yang, F. Wei, Z. Yang, Y. Su, S. Hu, Y. Chen, C.-M. Chan, W. Chen, et al., Parameter-efficient fine-tuning of large-scale pre-trained language models, *Nature Machine Intelligence* 5 (3) (2023) 220–235.
- [44] X. L. Li, P. Liang, Prefix-tuning: Optimizing continuous prompts for generation, in: Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing, 2021, pp. 4582–4597.
- [45] E. J. Hu, Y. Shen, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang, W. Chen, Lora: Low-rank adaptation of large language models, in: International Conference on Learning Representations, 2021.
URL <https://arxiv.org/abs/2106.09685>
- [46] N. Houlsby, A. Giurgiu, S. Jastrzebski, B. Morrone, Q. De Laroussilhe, A. Gesmundo, M. Attariyan, S. Gelly, Parameter-efficient transfer learning for NLP, in: International Conference on Machine Learning, 2019. [arXiv:1902.00751](https://arxiv.org/abs/1902.00751).
URL <https://arxiv.org/abs/1902.00751>
- [47] J. Hu, X. Xu, Z. Zou, LoRA-MedSAM: Efficient medical image segmentation, in: International Conference on Medical Imaging and Computer-Aided Diagnosis, Springer, 2024, pp. 154–164.
- [48] P. Lu, B. Peng, H. Cheng, M. Galley, K.-W. Chang, Y. N. Wu, S.-C. Zhu, J. Gao, Chameleon: Plug-and-play compositional reasoning with large language models, *Advances in Neural Information Processing Systems* 36 (2023) 43447–43478.
- [49] S. Woo, J. Park, J.-Y. Lee, I. S. Kweon, Cbam: Convolutional block attention module, in: Proceedings of the European conference on computer vision (ECCV), 2018, pp. 3–19.
- [50] D. Gutman, N. C. Codella, E. Celebi, B. Helba, M. Marchetti, N. Mishra, A. Halpern, Skin lesion analysis toward melanoma detection

tion: A challenge at the international symposium on biomedical imaging (ISBI) 2016, hosted by the international skin imaging collaboration (ISIC), arXiv preprint arXiv:1605.01397 (2016).

URL <https://arxiv.org/abs/1605.01397>

- [51] N. C. Codella, D. Gutman, M. E. Celebi, B. Helba, M. A. Marchetti, S. W. Dusza, A. Kalloo, K. Liopyris, N. Mishra, H. Kittler, et al., Skin lesion analysis toward melanoma detection: A challenge at the 2017 international symposium on biomedical imaging ISBI, hosted by the international skin imaging collaboration (ISIC), in: 2018 IEEE 15th International Symposium on Biomedical Imaging (ISBI), 2018, pp. 168–172. doi:10.1109/ISBI.2018.8363547.
- [52] N. Codella, V. Rotemberg, P. Tschandl, M. E. Celebi, S. Dusza, D. Gutman, B. Helba, A. Kalloo, K. Liopyris, M. Marchetti, H. Kittler, A. Halpern, Skin lesion analysis toward melanoma detection 2018: A challenge hosted by the international skin imaging collaboration (ISBI), arXiv preprint arXiv:1902.03368 (2019).
URL <https://arxiv.org/abs/1902.03368>
- [53] Y. Ji, H. Bai, C. Ge, J. Yang, Y. Zhu, R. Zhang, Z. Li, L. Zhanng, W. Ma, X. Wan, et al., Amos: A large-scale abdominal multi-organ benchmark for versatile medical image segmentation, Advances in neural information processing systems 35 (2022) 36722–36732.