

# Interactive Medical Image Segmentation: A Benchmark Dataset and Baseline

Junlong Cheng<sup>1,2,¶</sup> Bin Fu<sup>1</sup> Jin Ye<sup>1,3</sup> Guoan Wang<sup>1,4</sup> Tianbin Li<sup>1</sup>  
 Haoyu Wang<sup>1,5</sup> Ruoyu Li<sup>2</sup> He Yao<sup>2</sup> Junren Chen<sup>2</sup> JingWen Li<sup>6</sup>  
 Yanzhou Su<sup>1</sup> Min Zhu<sup>2,§</sup> Junjun He<sup>1,‡</sup>

<sup>1</sup>Shanghai AI Laboratory, General Medical Artificial Intelligence

<sup>2</sup>Sichuan University, School of Computer Science

<sup>3</sup>Monash University

<sup>4</sup>East China Normal University, School of computer science and technology

<sup>5</sup>Shanghai Jiao Tong University, School of biomedical engineering

<sup>6</sup>Xinjiang University, School of Computer Science and Technology

<sup>¶</sup>Main technical contribution    <sup>‡</sup>Corresponding authors    <sup>§</sup>Project lead  
 {chengjunlong, yejin, hejunjun, litianbin}@pjlab.org.cn

## Abstract

Interactive Medical Image Segmentation (IMIS) has long been constrained by the limited availability of large-scale, diverse, and densely annotated datasets, which hinders model generalization and consistent evaluation across different models. In this paper, we introduce the IMed-361M benchmark dataset, a significant advancement in general IMIS research. First, we collect and standardize over 6.4 million medical images and their corresponding ground truth masks from multiple data sources. Then, leveraging the strong object recognition capabilities of a vision foundational model, we automatically generated dense interactive masks for each image and ensured their quality through rigorous quality control and granularity management. Unlike previous datasets, which are limited by specific modalities or sparse annotations, IMed-361M spans 14 modalities and 204 segmentation targets, totaling 361 million masks—an average of 56 masks per image. Finally, we developed an IMIS baseline network on this dataset that supports high-quality mask generation through interactive inputs, including clicks, bounding boxes, text prompts, and their combinations. We evaluate its performance on medical image segmentation tasks from multiple perspectives, demonstrating superior accuracy and scalability compared to existing interactive segmentation models. To facilitate research on foundational models in medical computer vision, we release the IMed-361M and model at <https://github.com/uni-medical/IMIS-Bench>.

## 1 Introduction

Interactive Medical Image Segmentation (IMIS) enables clinicians or users to guide the model by marking points, lines, or regions on an image, resulting in segmentation outcomes that better meet clinical needs [1, 2]. This user-informed segmentation approach not only optimizes results throughout the segmentation process, aiding clinicians in precisely localizing targets during diagnosis and treatment [5, 6, 7], but also addresses the limitations of fully automated segmentation models in generalizing to unseen object classes (i.e., zero-shot learning) [3, 4, 8, 9, 10, 11], providing clinicians with a flexible adjustment mechanism that significantly enhances segmentation reliability. However, the advancement of IMIS advancement faces a critical bottleneck: the lack of high-quality, diverse, and large-scale datasets that capture the complex clinical environments [12, 13]. Existing public datasets (such as Totalsegmentator [14], AutoPET [27], and CoNIC [16]) are often tailored to specific

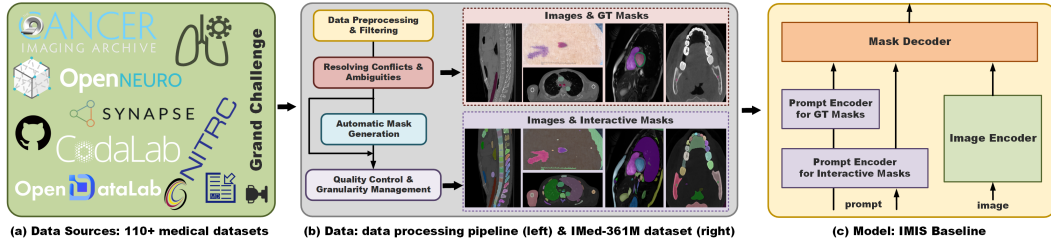


Figure 1: We collected 110 medical image datasets from various sources and generated the IMed-361M dataset, which contains over 361 million masks, through a rigorous and standardized data processing pipeline. Using this dataset, we developed the IMIS baseline network.

modalities or tasks, thus constraining the generalization capabilities of models across diverse medical scenarios.

Recently, as a foundational model for general visual segmentation, the Segment Anything Model (SAM) has garnered widespread attention [17, 18]. SAM integrates simple user interactions (such as points or bounding boxes) into the model’s learning process and leverages pre-training on large-scale datasets to achieve cross-domain and multi-task transferability [19, 20, 21]. Although SAM has been pre-trained on over 1 billion natural image masks, its direct application in medical imaging remains limited due to the substantial domain differences between natural and medical images, as well as the lack of medical imaging-specific knowledge. To address these issues, current methods such as MedSAM [22] and SAM-Med2D [24] attempt to expand dataset size and fine-tune SAM to create IMIS models tailored for clinical scenarios. While these approaches have shown some progress, their ability to ‘segment anything’ in medical imaging remains limited. This limitation is mainly due to the lack of densely masks in existing medical datasets compared to the SA-1B dataset [17] used to train SAM [22, 24, 25]. For instance, the COSMOS dataset [25] contains an average of only 5.7 masks per image, which restricting the model’s capacity for dense segmentation and hindering comprehensive, fine-grained interaction. Moreover, many methods are evaluated only on specific modalities or with limited interaction strategies, further restricting the comprehensiveness and reliability of IMIS model evaluation [22, 23, 25, 28].

Addressing these limitations requires the development of a high-quality IMIS benchmark dataset, which is essential for advancing foundational models in medical imaging [6, 29, 30, 31, 32]. An ideal IMIS benchmark dataset should meet three core criteria: **(1) Large-scale.** The dataset should be large enough to fully support deep learning model training, enabling the model to effectively capture medical features; **(2) Diversity.** The dataset should encompass various medical imaging modalities and complex clinical scenarios to ensure the model’s ability to generalize across modalities and tasks; **(3) High-quality and densely masks.** Accurate segmentation of complex medical images relies on high-quality annotations, and densely masks further enhance the model’s capability to “segment anything” in medical imaging. Unfortunately, existing public medical segmentation datasets [22, 24, 25, 47] do not fully meet these standards, limiting their ability to comprehensively support and evaluate IMIS models.

In this work, we introduce IMed-361M, a benchmark dataset specifically designed for IMIS tasks. As shown in Fig. 1, this dataset integrates public and private data sources and utilizes foundational models [17] for automated annotation to generating dense masks for each image. A standardized data processing workflow is ensure the high quality and consistency of all masks. IMed-361M achieves unprecedented scale, diversity, and mask quality, comprising 6.4 million images spanning 14 imaging modalities and 204 targets, with a total of 361 million masks, averaging 56 masks per image. IMed-361M effectively solves the problems of small data size and sparse annotations, providing data support for IMIS model training.

Additionally, we develop an IMIS baseline model and conduct a comprehensive evaluation of its performance across various medical scenarios, including its effectiveness on different modalities, anatomical structures, and organs, as well as the impact of various interaction strategies on model outcomes. This analysis provided an in-depth understanding of the strengths and limitations of different interactive segmentation methods, establishing a fair and consistent framework for evaluating IMIS model performance. We anticipate that the IMed-361M dataset and baseline model will drive

the widespread adoption of IMIS technology in clinical practice, accelerating the healthcare industry’s transition toward intelligence and automation.

## 2 Related Work

This section provides an overview of datasets used for various medical image segmentation tasks and reviews the advancements in IMIS algorithms.

**Datasets.** Due to variations in imaging protocols, medical images often have multidimensional characteristics (e.g., 2D, 3D), making pixel-level segmentation labeling for specific organs or lesions both expertise-dependent and time-consuming. Early work focused on single-organ or single-lesion annotated datasets [14, 27, 56, 57, 32, 59, 8]. These single-task datasets provide foundational data for training automatic segmentation models, significantly enhancing model performance in segmenting specific organs or lesions. With the growth of multi-task requirements, multi-organ and multi-tissue annotated datasets [14, 55, 61, 62] have emerged, facilitating the exploration of deep learning models in multi-region segmentation tasks [63, 64]. In recent years, visual foundational models have demonstrated impressive cross-domain and multi-task transfer capabilities, motivating researchers to build large-scale, multimodal datasets for medical foundational models. The most common approach is to integrate data from various sources [22, 24, 25, 47] to bridge the scale gap compared to natural image datasets. However, due to the sparsity of annotations in the source data, the combined dataset still lags behind natural image datasets in terms of mask quantity and density. For example, the SA-1B [17] dataset averages 100 masks per image, whereas large-scale medical datasets like COSMOS [25] and SA-Med2D-20M [39] have fewer than 5 masks per image on average. This gap limits the ability of medical datasets to support dense segmentation tasks and fine-grained interactions, as well as the applicability of foundational models to “segment anything” in the medical imaging domain.

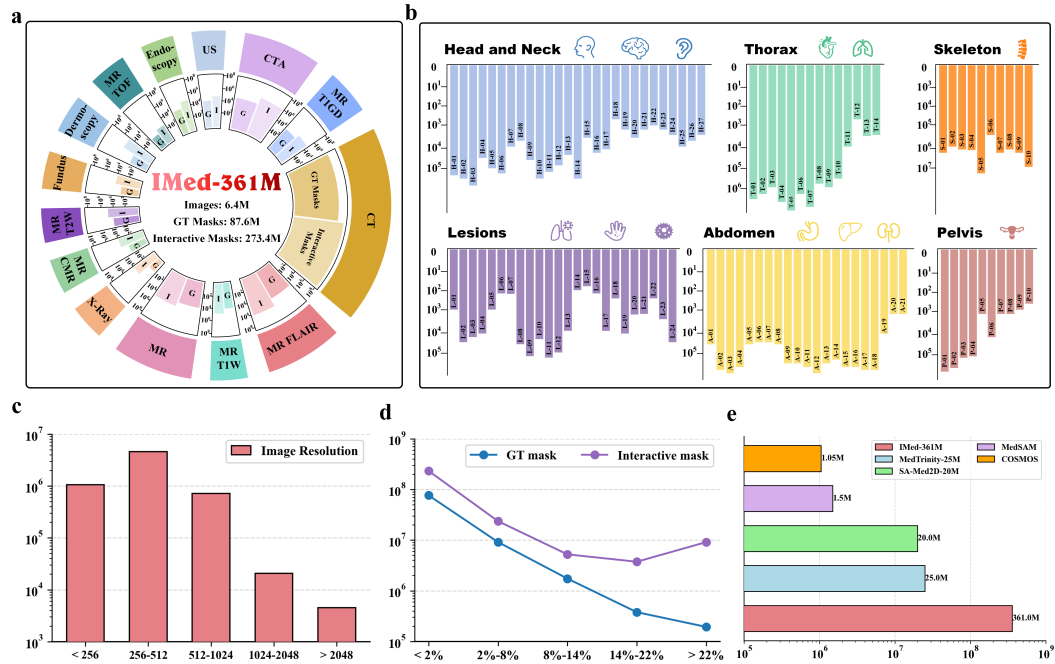


Figure 2: Overview of the IMed-361M dataset. (a) Number of images and masks for each modality. (b) Information on six anatomical structures. (c) Distribution of image resolutions. (d) Analysis of mask proportions. (e) Comparison with other existing public datasets.

**Algorithms.** IMIS methods are generally more effective than automated methods at generating high-quality results that meet clinical requirements, especially in scenarios involving diverse imaging protocols, complex pathological variations, and ambiguous lesion boundaries [70]. User interactions can take several typical forms, such as scribbles [71, 72], bounding boxes [17, 77], clicks [2, 73], or

language prompts [17, 74, 75]. Traditional methods approach this task through energy by minimizing energy on a regular pixel grid, capturing low-level appearance features with unary potentials and promoting consistent segmentation outputs with pairwise or higher-order potentials [72, 76]. With advancements in deep learning, researchers have explored using user prompts directly as network input features, often in combination with 2D or 3D neural networks to produce segmentation results [1, 2, 17, 71, 72, 73, 77]. Leveraging the advantage of pretraining on large-scale datasets, the Segment Anything Model (SAM) has become a benchmark for IMIS. For example, SAM can serve as a powerful pretrained encoder-decoder model fine-tuned for specific tasks [27, 65]. However, these approaches essentially revert to the design of fully automated segmentation models. Another category of methods retains the interactive aspect of segmentation, focusing on parameter-efficient fine-tuning (PEFT) techniques [66, 67]. However, these methods consider only a limited of interaction strategies, and inconsistencies in evaluation have affected the comparability and reliability of their results. To achieve robust segmentation capabilities across diverse medical imaging datasets, some approaches fine-tune SAM’s decoder using large-scale medical datasets [22, 24, 25]. Although these methods demonstrate significant advantages across different medical modalities and tasks, they have not yet been consistently benchmarked against similar methods and lack a thorough examination of how different interaction strategies affect outcomes.

In this work, we evaluate these state-of-the-art methods on the IMed-361M dataset, providing a comprehensive and fair comparison, while also exploring the impact of different interaction strategies on medical image segmentation.

### 3 IMIS Benchmark Dataset: IMed-361M

We present IMed-361M, the first large-scale IMIS dataset, that significantly surpasses existing datasets in terms of scale, diversity, quality, and density. This section details the data collection and pre-processing procedures, along with an in-depth analysis of the dataset’s advantages and potential applications.

#### 3.1 Data Collection and Pre-processing

**Data collection.** We integrate over 110 publicly available medical image segmentation datasets from globally recognized platforms such as TCIA<sup>1</sup>, OpenNeuro<sup>2</sup>, NITRC<sup>3</sup>, Grand Challenge<sup>4</sup>, Synapse<sup>5</sup>, CodaLab<sup>6</sup>, and GitHub<sup>7</sup>, covering both 2D and 3D images and variety of formats (e.g., .jpg, .npy, .nii). Additionally, we collaborate with several medical institutions, acquiring diverse clinical data through rigorous ethical review processes, which further enriches the dataset. A detailed list of data sources is provided in the supplementary material.

**Preprocessing and filtering.** We first standardize all collected medical images according to the SA-Med2D-20M [39] protocol, converting the corresponding ground truth (GT) into one-hot encoding and storing it in .npz files in compressed sparse row (CSR) format [40]. Then, we apply the following exclusion criteria: (1) exclude 3D slice images and their corresponding masks with an aspect ratio greater than 1.5 to avoid interference from images with excessive geometric distortion during model training; (2) exclude masks where the foreground area accounts for less than one-thousandth of the total pixels, as such small foreground regions are prone to being lost during resizing, making it challenging to generate effective prompts.

**Resolving conflicts and ambiguities.** To address conflicts and ambiguities in the GT, we first standardize expressions for the same target. For instance, both "lung nodule" and "pulmonary nodule" are renamed to "lung nodule" for consistency. We then manually review and correct misalignments and informational errors in the dataset to ensure annotation accuracy. Finally, for annotations with

<sup>1</sup><https://www.cancerimagingarchive.net>

<sup>2</sup><https://openneuro.org>

<sup>3</sup><https://www.nitrc.org>

<sup>4</sup><https://grand-challenge.org>

<sup>5</sup><https://www.synapse.org>

<sup>6</sup><https://codalab.org>

<sup>7</sup><https://github.com>



multiple connected components, we differentiate and label them based on clinical needs to avoid potential misunderstandings that could arise from single-point interactions.

Through this process, we collected 6.4 million images and 87.6 million GT masks that were manually annotated and corrected. Although this already exceeds the data volume of most existing segmentation datasets, it still does not meet our fundamental requirement for the IMIS dataset: diverse and densely masks. Therefore, we further introduce interactive masks to address the above issues.

### 3.2 Interactive Masks

**Automatic mask generation.** The automatic mask generation method employed by SAM has been proven to be high-quality and effective [17, 43, 44]. We leverage SAM’s object-awareness capability to generate as many masks as possible for each image. Specifically, we employ a 32×32 point grid to guide the model, generating a set of candidate masks for each point corresponding to potential objects of interest. The generated masks are refined and optimized using the following strategies: (1) Confidence filtering: Only masks with an Intersection over Union (IoU) prediction confidence score above 0.85 are retained to ensure segmentation accuracy. (2) Non-Maximum Suppression (NMS) [45]: For overlapping masks with an IoU greater than 0.7, only the mask with the highest confidence score is kept to eliminate redundancy. (3) Remove background masks: Masks covering more than 80% of the area are discarded, as such large-coverage foreground masks are almost nonexistent in the GT.

**Quality control and granularity management.** Based on our observations, the generated masks often fail to fully separate structures with unclear boundaries, such as the left atrium, myocardium, and right atrium in the heart. Additionally, masks for dispersed structures like the intestines are often identified as multiple separate objects. To address these issues, we utilize the original GT (described in Section 3.1) to correct the generated masks. We focus on two situations: (1) If the GT contains multi-connected regions, we directly replace the corresponding regions in the generated mask with the multi-connected regions from the GT; (2) We iterate over all single-connected regions in the GT, generating a minimum bounding box for each region. If there is a region in the generated mask with a bounding box that overlaps more than 95% with this GT bounding box, we retain the generated mask region; otherwise, we replace the region with the GT mask. Finally, we apply morphological operations (e.g., erosion and dilation) to remove noise and fill small holes [46].

We ultimately obtained 273 million "interactive masks", which can be used to train interactive segmentation models and cover nearly all identifiable objects in medical images. While we acknowledge that some masks may lack direct clinical relevance, they provide considerable diversity and density in the identified regions, helping the model learn a variety of interactive operations and enhancing its generalization ability across various tasks.

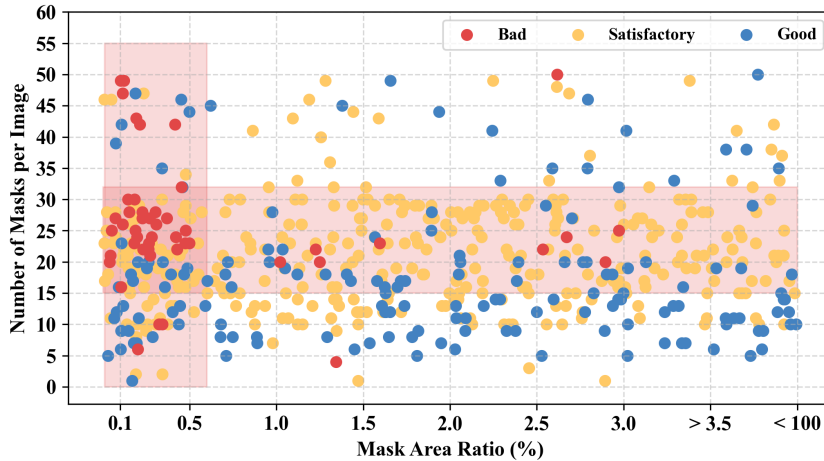


Figure 3: Evaluation of the quality of interactive masks.

### 3.3 Statistics and Analysis

**Data scale.** As shown in Fig.2 (a), the IMed-361M dataset contains 6.4 million images, 87.6 million GT, and 273.4 million interactive masks, averaging 56 masks per image. This makes it the largest publicly available, multimodal, interactive medical image segmentation dataset to date. Compared to the MedTrinity-25M [47] dataset, our dataset provides 14.4 times the number of masks (Fig.2 (e)).

**Diversity.** The IMed-361M dataset covers 14 imaging modalities and 204 segmentation targets, including organs and lesions, as shown in Fig.2 (b). We categorize the GT into six groups: Head and Neck, Thorax, Skeleton, Abdomen, Pelvis, and Lesions, covering nearly all parts of the human body. For clarity, we merge certain positional information when appropriate; for example, “left lung” and “right lung” are combined into “lung.” Detailed category information can be found in the supplementary materials. Fig.2 (c) and (d) show the distribution of image resolutions and mask coverage within the dataset. Over 83% of the images have resolutions between  $256 \times 256$  and  $1024 \times 1024$ , ensuring broad applicability and compatibility across various research scenarios. Additionally, most masks occupy less than 2% of the image area, reflecting the typically fine granularity of medical segmentation. Our interactive masks provide over one million instances across different coverage intervals, significantly enhancing the diversity and density of the dataset.

**Mask quality.** We randomly select five images from each subset of IMed-361M (a total of 550 images) and invited four radiologists to evaluate the quality of the corresponding interactive masks. We categorize the mask quality into three levels: (1) "Good": the mask closely matches the manual annotations; (2) "Satisfactory": the mask has some visible errors (e.g., incomplete contours) but is generally acceptable; (3) "Bad": the mask has significant errors (e.g., over-segmentation or under-segmentation) and is unsuitable for clinical segmentation tasks. The evaluation results show that 156 images are rated as "Good," 345 as "Satisfactory," and 49 as "Bad". Statistical analysis reveals that the masks rated as "Bad" mainly come from 18 datasets, concentrated within the red-boxed area in Fig.3. Radiologists note that some excessively annotated structures hold little clinical relevance. Additionally, some masks include elements unrelated to human organs or lesions, such as letters in X-ray backgrounds or hair in dermatoscopic images. Based on this analysis, we manually remove the excessively annotated masks from these 18 datasets and apply a 0.5% foreground filter rate to retain only the final valid masks. Finally, we retain those masks unrelated to human organs or lesions to enhance the model’s adaptability in different scenarios.

## 4 IMIS Baseline Network

### 4.1 Model Design

We adopt a strategy similar to fine-tuning SAM [22, 24, 25] to establish the IMIS baseline, providing a performance benchmark for future work. As shown in Fig.4, IMIS-Net has three key components: an image encoder for extracting image features, a prompt encoder for integrating user interaction information, and a mask decoder for generating segmentation results using image and prompt embeddings.

We select ViT-base [48] as the image encoder. While larger ViT models (e.g., ViT-Large, ViT-Huge) [22, 25] offer slight accuracy improvements, they significantly increase computational costs, making them impractical for clinical applications. Additionally, most open-source fine-tuning methods also use ViT-base, ensuring fairness in our baseline. The image input size is  $1024 \times 1024 \times 3$ , with a patch size of  $16 \times 16 \times 3$ . For the prompt encoder, we consider three prompt types: points, boxes, and text. Points and boxes are represented by the sum of positional encoding and learned embeddings, while text is encoded using CLIP’s text encoder [49]. Text prompts follow the template: "A segmentation area of a [category]," covering over 200 organ and lesion categories to establish a benchmark for future multimodal segmentation research. The mask decoder uses Transformer decoder blocks, mapping image and prompt embeddings to the mask. The mask resolution is increased to  $256 \times 256$  through transposed convolutions and matched to the input size via bilinear interpolation. Notably, expanding the encoder can improve performance without significantly increasing network training parameters, thereby providing a simple solution to address the issue of model performance saturation.

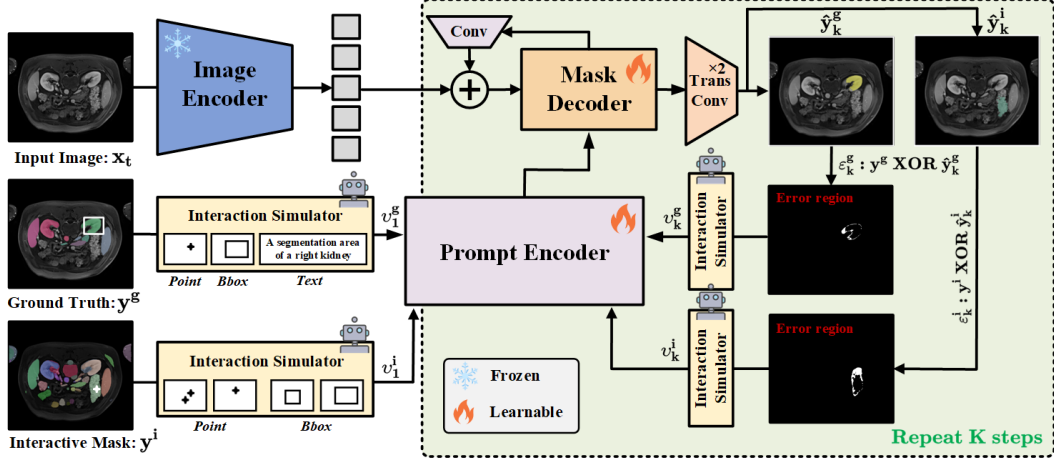


Figure 4: The training process of IMIS-Net simulates  $K$  consecutive steps of interactive segmentation.

## 4.2 Experiment Configuration

### Training strategy.

An overview is provided of the simulated continuous interactive segmentation training [22, 24, 50] in Fig.4. For a given segmentation task and medical image  $x_t$ , we first simulate a set of initial interactions  $u_1^g$  and  $u_1^i$  based on the corresponding ground truth  $y^g$  and interactive mask  $y^i$ , which include clicks, bboxes, and text input. The click points are uniformly sampled from the foreground regions of  $y^g$  or  $y^i$ , while the bboxes are defined as the smallest bounding box around the target, with an offset of 5 pixels added to each coordinate to simulate slight user bias during the interaction process. The entire training process involves  $K$  interactive training iterations (with  $K = 8$  in this paper). The model’s initial predictions are  $\hat{y}_1^g$  and  $\hat{y}_1^i$ . After the first prediction, we simulate subsequent corrections based on the previous predictions  $\hat{y}_k^g$  and  $\hat{y}_k^i$ , as well as the error region  $\varepsilon_k$  between the  $y^g$  and  $y^i$ , where  $k \in \{1, \dots, K\}$ . Additionally, we provide the low-resolution predicted mask from the previous prediction as an extra cue to the model. As can be seen, the image encoder only needs to encode the image once during the training, and subsequent interactive training only updates the prompt encoder and mask decoder parameters.

**Training setting.** We use the Adam optimizer [51] with a learning rate of  $2 \times 10^{-5}$ . The training was conducted on 72 NVIDIA 4090 GPUs with a batch size of 2. For each image, we randomly select 5 targets from the corresponding GT and interactive masks as supervision targets (if fewer than 5 targets are present, they are selected repeatedly). The image resolution was uniformly resized to  $1024 \times 1024$ , and pixel intensities are randomly scaled and shifted with a probability of 20%, where the scaling factor and offset are set to 0.2. We train the model for a total of 12 epochs on the IMed-361M dataset and select the final checkpoint as the model weights. A linear combination of Focal loss [52] and Dice loss [53] is used as the loss function, balancing their influence in a 20:1 ratio. Finally, the Dice score is used as the primary evaluation metric for this study.

### 4.3 Evaluation Model and Data

**Model.** We evaluate five visual foundation models [17, 18, 22, 24, 25], selected based on their pretraining on large-scale datasets and the availability of source code. Among them, MedSAM, SAM-Med2D, and Huang et al. [25] are specifically designed for IMIS, while SAM [17] and SAM-2 [18] are pre-trained on massive natural image and video datasets. To ensure fairness in the evaluation, all models used ViT-Base as the encoder. It is important to note that, MedSAM and Huang et al. [25] only support bbox-based interactivity (using the models released by the authors), while the other models support points, bboxes, and combinations of both.

**Data.** The data used to evaluate the model is divided into two categories: the test set and the external dataset. The test set contains 58,411 images, with a distribution consistent with that of IMed-361M. The external dataset includes data from three different modalities and clinical scenarios: SegThor [54],

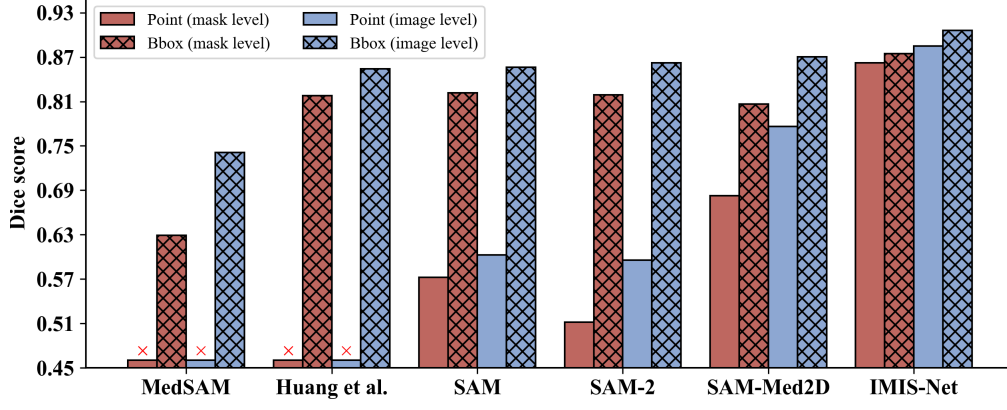


Figure 5: Comparison of IMIS-Net with existing foundation models, with performance statistics at both image and mask levels.

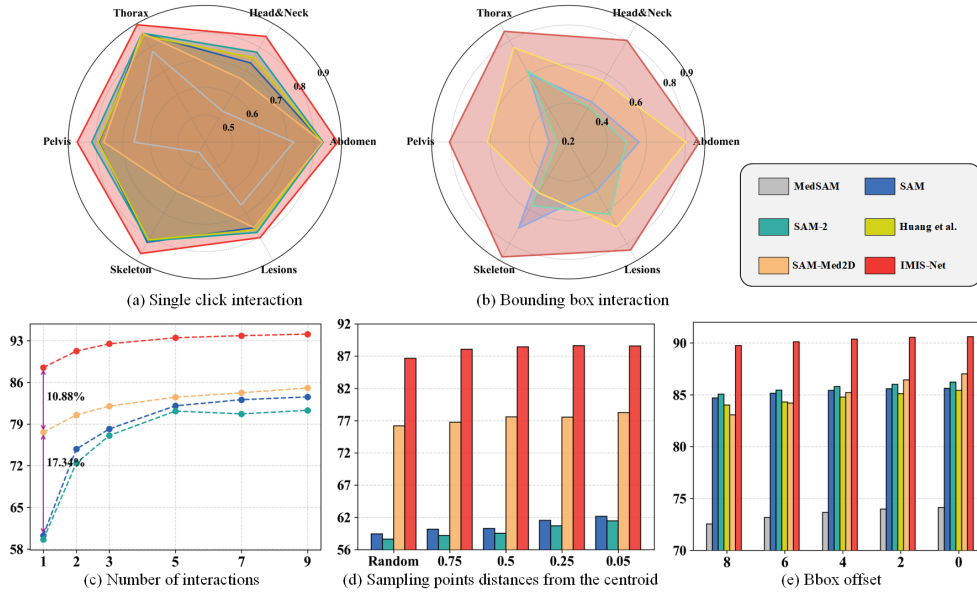


Figure 6: Comparison of segmentation performance across different anatomical structures under single-click (a) and bounding box (b) interactions. (c) Changes in segmentation performance with increasing interaction numbers. (d) and (e) Impact of click position and bounding box offset on performance.

which is used for segmentation of thoracic organs at risk in 516 CT images; TotalSegmentator MRI [55], which is used for segmenting organs in 2,147 MRI images; and ISLES<sup>8</sup>, which belongs to the MRI FLAIR modality and includes 364 images of ischemic stroke lesions. These datasets were not used during the training of IMIS-Net.

## 5 Results

### 5.1 Main Results

Fig.5 shows the performance evaluation of IMIS-Net compared to other vision foundation models on the single-interaction segmentation task. During testing, all models used the same points and bounding boxes as prompts. The results indicate that IMIS-Net outperforms other models in both

<sup>8</sup><https://www.isles-challenge.org/>

Dataset	Category	SAM	SAM 2	MedSAM	Huang et al.	SAM-Med2D	IMIS-Net
ISLES	Ischemic Stroke Lesion	55.92	60.14	59.90	56.79	68.22	<b>71.78</b>
SegThor	Esophagus	78.19	86.60	47.17	76.48	82.04	<b>89.17</b>
	Heart	91.41	92.34	81.16	92.02	<b>94.18</b>	81.27
	Aorta	82.36	83.69	55.68	84.57	79.12	<b>94.92</b>
	Trachea	93.76	94.24	72.40	<b>94.60</b>	85.10	90.04
	<b>Average</b>	84.46	85.86	60.55	84.73	86.43	<b>89.27</b>
TotalSegmentator MRI	Adrenal gland	77.19	<b>82.25</b>	56.58	68.52	66.35	77.53
	Aorta	76.88	79.90	55.98	80.08	74.91	<b>85.62</b>
	Colon	56.69	56.00	42.53	55.41	51.44	<b>73.47</b>
	Duodenum	69.30	74.81	55.64	72.38	69.67	<b>75.97</b>
	Gallbladder	81.33	84.82	68.15	83.46	<b>87.63</b>	82.93
	Iliopsoas	59.57	71.96	58.48	63.67	60.76	<b>72.55</b>
	Kidney	79.43	81.59	58.41	81.09	74.04	<b>82.72</b>
	Inferior vena cava	84.94	87.03	57.73	85.88	78.34	<b>87.35</b>
	Liver	86.93	88.70	68.14	87.35	86.75	<b>90.62</b>
	Pancreas	66.95	69.86	46.10	66.30	66.84	<b>70.84</b>
	Spleen	81.12	83.05	66.61	83.87	84.95	<b>86.74</b>
	Stomach	77.10	81.23	62.75	76.52	79.31	<b>81.45</b>
	<b>Average</b>	75.45	77.62	59.52	75.47	75.92	<b>79.06</b>

Table 1: Quantitative comparison results of IMIS-Net against five other interactive segmentation methods on three external datasets.

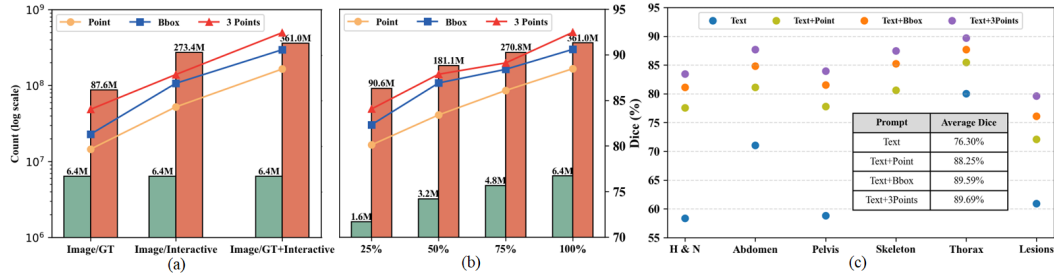


Figure 7: (a) and (b) segmentation performance increases with the increasing mask density or data scales up. (c) effectiveness of text and its combination with other prompts.

image and mask-level statistics. The bbox interaction of different models always outperforms the click interaction because the bbox can provide more boundary information. Notably, despite being pretrained on large-scale medical image datasets, MedSAM and SAM-Med2D still exhibit significant performance differences. This disparity is closely related to the scale and diversity of the pretraining datasets. As shown in Fig.6(a) and (b), the pre-training dataset for SAM-Med2D lacks samples of skeletal structures, resulting in poorer segmentation performance for these anatomical structures. Additionally, under the single-point prompt condition, SAM and SAM-2 achieve only Dice scores of 60.26% and 59.57%, respectively, likely due to the absence of medical knowledge in the pretraining data and limited interactive information constraining model performance.

We increase the number of interactions from 1 to 9 to observe the performance changes of the model (Fig.6 (c)). As expected, performance improves with more interactions, and the gap between models narrows. This is because multiple interactions reduce the task difficulty. We also examine the impact of click position and bbox offset on performance. As shown in Fig.6 (d) and (f), when the prompt point is closer to the centroid, the performance improvement is more significant, with SAM-2's Dice score increasing by 2.84%. Additionally, with bbox offset, all methods experience a performance decline of 0.85%-3.94%. Our model exhibited the smallest performance drop, making it more robust for practical use. In summary, leveraging the rich data and mask diversity of the IMed-361M dataset, our model delivers the best performance across various medical scenarios and interaction strategies. More experimental results and analysis can be found in the supplementary material.

## 5.2 External Dataset Evaluation

Tab.1 presents evaluation results on external datasets. IMIS-Net achieves the best average performance across datasets from three different sources and tasks. We highlight the Dice scores of the 12 major abdominal organs in the TotalSegmentator MRI dataset, where IMIS-Net outperforms in 10. On the



ISLES dataset, focusing on ischemic stroke lesion segmentation, our model surpasses the second-best SAM-Med2D by 3.56%. These results demonstrate the strong generalization ability of our method.

Decoder dimension	Image resolution	Prompt		Training parameters
		Point	Bbox	
768	256×256	0.8214	0.8469	29.68 M
768	512×512	0.8673	0.8968	29.68 M
256	1024×1024	0.8366	0.8497	5.52 M
512	1024×1024	0.8563	0.8729	15.19 M
768	1024×1024	0.8848	0.9060	29.68 M

Table 2: Ablation study of model design, including decoder dimension and image resolution.

### 5.3 Ablation Study

**Scaling Up Training Data.** As shown in Fig.4 (a) and (b), when trained only on the GT from IMed-361M, IMIS-Net performs poorly. However, with the addition of interactive masks, the Dice score increases rapidly. Similar performance trends are observed when training on datasets of different sizes, indicating that our method is scalable and performs better with more available training data.

**Text and Combination Prompts.** IMIS-Net, as a versatile interactive model, can utilize text prompts to achieve the segmentation of over 200 organs and lesions. As shown in Fig.4 (c), the model achieves a segmentation performance of 76.30% when using only text prompts. When both text and point prompts are combined, the average Dice score increases by 11.95%. Furthermore, after three rounds of click-based correction, the Dice score reaches 89.69%. These results show that combining different prompts synergistically enhances the model’s segmentation capabilities.

**Model Design.** We conduct an ablation study on different configurations of the model, including the impact of input resolution and decoder dimension on model performance. As shown in Tab.2, the model’s segmentation performance improves as the training image resolution increases. This higher input resolution allows for clearer visualization of lesions and organs in medical images. Additionally, when the decoder dimension is increased from 256 to 768, the model’s performance improves from 84.97% to 90.60%, with only a 24.16M increase in trainable parameters. These results demonstrate the scalability of IMIS-Net, showing that even when the model’s performance approaches saturation on larger datasets, further performance improvements can still be achieved by expanding the decoder.

## 6 Conclusion

In this work, we introduce IMed-361M, a benchmark dataset dedicated to interactive medical image segmentation. It includes a vast array of medical images across various modalities, extensive segmentation scenarios, and densely masks, surpassing all existing datasets that are limited to single tasks or simple integrations. Leveraging this data resource, we developed a general IMIS baseline model that enables users to generate segmentation results tailored to clinical needs through interactive methods, including clicks, bounding boxes, text prompts, and their combinations. We performed a comprehensive comparison of this baseline with existing foundational models, demonstrating that our model provides significant performance advantages and exhibits strong transferability in previously unseen scenarios. Notably, our approach typically requires fewer interactions to achieve comparable performance, enhancing its practicality in real-world applications.

The IMed-361M dataset will strongly facilitate the development of foundational models in medical imaging and lay the foundation for fair evaluation across different models. IMIS-Net provides general technical support for various clinical applications, accelerating the widespread application of AI technology in the medical field. Despite these achievements, we recognize that this work still faces several challenges. For instance, effectively obtaining semantic information for interactive masks and extending this approach to more comprehensive and finer-grained medical image analysis scenarios are areas that require further exploration and improvement in the future.

## References

- [1] Zhou T, Li L, Bredell G, et al. Volumetric memory network for interactive medical image segmentation[J]. *Medical Image Analysis*, 2023, 83: 102599.
- [2] Zhang J, Shi Y, Sun J, et al. Interactive medical image segmentation via a point-based interaction[J]. *Artificial Intelligence in Medicine*, 2021, 111: 101998.
- [3] Azad R, Aghdam E K, Rauland A, et al. Medical image segmentation review: The success of u-net[J]. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024.
- [4] Cheng J, Tian S, Yu L, et al. ResGANet: Residual group attention network for medical image classification and segmentation[J]. *Medical Image Analysis*, 2022, 76: 102313.
- [5] Chen F, Han H, Wan P, et al. Joint segmentation and differential diagnosis of thyroid nodule in contrast-enhanced ultrasound images[J]. *IEEE Transactions on Biomedical Engineering*, 2023, 70(9): 2722-2732.
- [6] Marinov Z, Jäger P F, Egger J, et al. Deep interactive segmentation of medical images: A systematic review and taxonomy[J]. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024.
- [7] Conze P H, Andrade-Miranda G, Singh V K, et al. Current and emerging trends in medical image segmentation with deep learning[J]. *IEEE Transactions on Radiation and Plasma Medical Sciences*, 2023, 7(6): 545-569.
- [8] Ronneberger O, Fischer P, Brox T. U-net: Convolutional networks for biomedical image segmentation[C]. *Medical Image Computing and Computer-Assisted Intervention–MICCAI 2015*, 2015: 234-241.
- [9] Chen J, Lu Y, Yu Q, et al. Transunet: Transformers make strong encoders for medical image segmentation[J]. *arXiv preprint [arXiv:2102.04306](https://arxiv.org/abs/2102.04306)*, 2021.
- [10] Zhou Z, Rahman Siddiquee M M, Tajbakhsh N, et al. Unet++: A nested u-net architecture for medical image segmentation[C]. *Medical Image Computing and Computer-Assisted Intervention–MICCAI 2018*, 2018: 3-11.
- [11] Xing Z, Ye T, Yang Y, et al. Segmamba: Long-range sequential modeling mamba for 3d medical image segmentation[C]. *International Conference on Medical Image Computing and Computer-Assisted Intervention*, 2024: 578-588.
- [12] Cascella M, Montomoli J, Bellini V, et al. Evaluating the feasibility of ChatGPT in healthcare: an analysis of multiple clinical and research scenarios[J]. *Journal of Medical Systems*, 2023, 47(1): 33.
- [13] Yuan M, Xia Y, Dong H, et al. Devil is in the queries: advancing mask transformers for real-world medical image segmentation and out-of-distribution localization[C]. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023: 23879-23889.
- [14] Wasserthal J, Breit H C, Meyer M T, et al. TotalSegmentator: robust segmentation of 104 anatomic structures in CT images[J]. *Radiology: Artificial Intelligence*, 2023, 5(5).
- [15] Gatidis S, Hepp T, Früh M, et al. A whole-body FDG-PET/CT dataset with manually annotated tumor lesions[J]. *Scientific Data*, 2022, 9(1): 601.
- [16] Graham S, Jahanifar M, Vu Q D, et al. Conic: Colon nuclei identification and counting challenge 2022[J]. *arXiv preprint [arXiv:2111.14485](https://arxiv.org/abs/2111.14485)*, 2021.
- [17] Kirillov A, Mintun E, Ravi N, et al. Segment anything[C]. *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023: 4015-4026.
- [18] Ravi N, Gabeur V, Hu Y T, et al. Sam 2: Segment anything in images and videos[J]. *arXiv preprint [arXiv:2408.00714](https://arxiv.org/abs/2408.00714)*, 2024.

- [19] Ke L, Ye M, Danelljan M, et al. Segment anything in high quality[J]. Advances in Neural Information Processing Systems, 2024, 36.
- [20] Cen J, Zhou Z, Fang J, et al. Segment anything in 3D with nerfs[J]. Advances in Neural Information Processing Systems, 2023, 36: 25971-25990.
- [21] Wang D, Zhang J, Du B, et al. Samrs: Scaling-up remote sensing segmentation dataset with segment anything model[J]. Advances in Neural Information Processing Systems, 2024, 36.
- [22] Ma J, He Y, Li F, et al. Segment anything in medical images[J]. Nature Communications, 2024, 15(1): 654.
- [23] Zhang Y, Shen Z, Jiao R. Segment anything model for medical image segmentation: Current applications and future directions[J]. Computers in Biology and Medicine, 2024: 108238.
- [24] Cheng J, Ye J, Deng Z, et al. Sam-med2d[J]. arXiv preprint [arXiv:2308.16184](https://arxiv.org/abs/2308.16184), 2023.
- [25] Huang Y, Yang X, Liu L, et al. Segment anything model for medical images?[J]. Medical Image Analysis, 2024, 92: 103061.
- [26] Miao, Juzheng C, Cheng Z, et al. Cross prompting consistency with segment anything model for semi-supervised medical image segmentation[C]. International Conference on Medical Image Computing and Computer-Assisted Intervention, 2024: 167–177.
- [27] Hu M, Li Y, Yang X. SkinSAM: Adapting the segmentation anything model for skin cancer segmentation[C]. Medical Imaging 2024: Image Perception, Observer Performance, and Technology Assessment, SPIE, 2024, 12929: 174-184.
- [28] Mazurowski M A, Dong H, Gu H, et al. Segment anything model for medical image analysis: An experimental study[J]. Medical Image Analysis, 2023, 89: 102918.
- [29] Zhang S, Metaxas D. On the challenges and perspectives of foundation models for medical image analysis[J]. Medical Image Analysis, 2024, 91: 102996.
- [30] Prabhod K J, Gadhiraaju A. Foundation Models in Medical Imaging: Revolutionizing Diagnostic Accuracy and Efficiency[J]. Journal of Artificial Intelligence Research and Applications, 2024, 4(1): 471-511.
- [31] Chen R J, Ding T, Lu M Y, et al. Towards a general-purpose foundation model for computational pathology[J]. Nature Medicine, 2024, 30(3): 850-862.
- [32] Schäfer R, Nicke T, Höfener H, et al. Overcoming data scarcity in biomedical imaging with a foundational multi-task model[J]. Nature Computational Science, 2024, 4(7): 495-509.
- [33] Wang A, Liu A, Zhang R, et al. REVISE: A tool for measuring and mitigating bias in visual datasets[J]. International Journal of Computer Vision, 2022, 130(7): 1790-1810.
- [34] Albahri A S, Duham A M, Fadhel M A, et al. A systematic review of trustworthy and explainable artificial intelligence in healthcare: Assessment of quality, bias risk, and data fusion[J]. Information Fusion, 2023, 96: 156-191.
- [35] Khan M W, Sheng H, Zhang H, et al. RVD: A handheld device-based fundus video dataset for retinal vessel segmentation[J]. Advances in Neural Information Processing Systems, 2024, 36.
- [36] Jin K, Huang X, Zhou J, et al. Fives: A fundus image dataset for artificial intelligence-based vessel segmentation[J]. Scientific Data, 2022, 9(1): 475.
- [37] Spanhol F A, Oliveira L S, Petitjean C, et al. A dataset for breast cancer histopathological image classification[J]. IEEE Transactions on Biomedical Engineering, 2015, 63(7): 1455-1462.
- [38] Brancati N, Anniciello A M, Pati P, et al. Bracs: A dataset for breast carcinoma subtyping in H&E histology images[J]. Database, 2022: baac093.
- [39] Ye J, Cheng J, Chen J, et al. Sa-med2d-20m dataset: Segment anything in 2D medical imaging with 20 million masks[J]. arXiv preprint [arXiv:2311.11969](https://arxiv.org/abs/2311.11969), 2023.

- [40] Borštnik U, VandeVondele J, Weber V, et al. Sparse matrix multiplication: The distributed block-compressed sparse row library[J]. *Parallel Computing*, 2014, 40(5-6): 47-58.
- [41] Gong H, Huang W, Zhang H, et al. Intensity confusion matters: An intensity-distance guided loss for bronchus segmentation[C]. *2024 IEEE International Conference on Multimedia and Expo (ICME)*, 2024: 1-6.
- [42] Zeng L L, Gao K, Hu D, et al. SS-TBN: A semi-supervised tri-branch network for COVID-19 screening and lesion segmentation[J]. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023, 45(8): 10427-10442.
- [43] He C, Li K, Zhang Y, et al. Weakly-supervised concealed object segmentation with SAM-based pseudo labeling and multi-scale feature grouping[J]. *Advances in Neural Information Processing Systems*, 2024, 36.
- [44] Zhang Y, Zhou T, Wang S, et al. Input augmentation with SAM: Boosting medical image segmentation with segmentation foundation model[C]. *International Conference on Medical Image Computing and Computer-Assisted Intervention*, 2023: 129-139.
- [45] Symeonidis C, Mademlis I, Pitas I, et al. Neural attention-driven non-maximum suppression for person detection[J]. *IEEE Transactions on Image Processing*, 2023, 32: 2454-2467.
- [46] Tummala S K, et al. Morphological operations and histogram analysis of SEM images using Python[J]. *Indian Journal of Engineering and Materials Sciences (IJEMS)*, 2023, 29(6): 796-800.
- [47] Xie Y, Zhou C, Gao L, et al. Medtrinity-25m: A large-scale multimodal dataset with multigranular annotations for medicine[J]. *arXiv preprint [arXiv:2408.02900](https://arxiv.org/abs/2408.02900)*, 2024.
- [48] Dosovitskiy A. An image is worth 16x16 words: Transformers for image recognition at scale[J]. *arXiv preprint [arXiv:2010.11929](https://arxiv.org/abs/2010.11929)*, 2020.
- [49] Radford A, Kim J W, Hallacy C, et al. Learning transferable visual models from natural language supervision[C]. *International Conference on Machine Learning*, 2021: 8748-8763.
- [50] Sofiuk K, Petrov I A, Konushin A. Reviving iterative training with mask guidance for interactive segmentation[C]. *2022 IEEE International Conference on Image Processing (ICIP)*, 2022: 3141-3145.
- [51] Zhang Z. Improved adam optimizer for deep neural networks[C]. *2018 IEEE/ACM 26th International Symposium on Quality of Service (IWQoS)*, 2018: 1-2.
- [52] Ross T-Y, Dollár G. Focal loss for dense object detection[C]. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017: 2980-2988.
- [53] Milletari F, Navab N, Ahmadi S A. V-net: Fully convolutional neural networks for volumetric medical image segmentation[C]. *2016 Fourth International Conference on 3D Vision (3DV)*, 2016: 565-571.
- [54] Lambert Z, Petitjean C, Dubray B, et al. Segthor: Segmentation of thoracic organs at risk in CT images[C]. *2020 Tenth International Conference on Image Processing Theory, Tools and Applications (IPTA)*, 2020: 1-6.
- [55] D'Antonoli T A, Berger L K, Indrakanti A K, et al. TotalSegmentator MRI: Sequence-Independent Segmentation of 59 Anatomical Structures in MR images[J]. *arXiv preprint [arXiv:2405.19492](https://arxiv.org/abs/2405.19492)*, 2024.
- [56] Clark K, Vendt B, Smith K, et al. The Cancer Imaging Archive (TCIA): Maintaining and operating a public information repository[J]. *Journal of Digital Imaging*, 2013, 26: 1045-1057.
- [57] Pedrosa J, Aresta G, Ferreira C, et al. LNDb challenge on automatic lung cancer patient management[J]. *Medical Image Analysis*, 2021, 70: 102027.
- [58] Rudyanto R D, Kerkstra S, Van Rikxoort E M, et al. Comparing algorithms for automated vessel segmentation in computed tomography scans of the lung: the VESSEL12 study[J]. *Medical Image Analysis*, 2014, 18(7): 1217-1232.

- [59] Antonelli M, Reinke A, Bakas S, et al. The medical segmentation decathlon[J]. *Nature Communications*, 2022, 13(1): 4128.
- [60] Jha D, Smedsrud P H, Riegler M A, et al. Kvasir-seg: A segmented polyp dataset[C]. *MultiMedia Modeling: 26th International Conference, MMM 2020, Daejeon, South Korea, January 5-8, 2020, Proceedings, Part II*. Springer, 2020: 451-462.
- [61] Zingman I, Stierstorfer B, Lempp C, et al. Learning image representations for anomaly detection: application to discovery of histological alterations in drug development[J]. *Medical Image Analysis*, 2024, 92: 103067.
- [62] Podobnik G, Strojanić P, Peterlin P, et al. HaN-Seg: The head and neck organ-at-risk CT and MR segmentation dataset[J]. *Medical Physics*, 2023, 50(3): 1917-1927.
- [63] Isensee F, Jaeger P F, Kohl S A A, et al. nnU-Net: a self-configuring method for deep learning-based biomedical image segmentation[J]. *Nature Methods*, 2021, 18(2): 203-211.
- [64] Huang Z, Wang H, Deng Z, et al. Stu-net: Scalable and transferable medical image segmentation models empowered by large-scale supervised pre-training[J]. *arXiv preprint [arXiv:2304.06716](https://arxiv.org/abs/2304.06716)*, 2023.
- [65] Li Y, Hu M, Yang X. Polyp-sam: Transfer sam for polyp segmentation[C]. *Medical Imaging 2024: Computer-Aided Diagnosis, SPIE*, 2024, 12927: 759-765.
- [66] Wu J, Ji W, Liu Y, et al. Medical sam adapter: Adapting segment anything model for medical image segmentation[J]. *arXiv preprint [arXiv:2304.12620](https://arxiv.org/abs/2304.12620)*, 2023.
- [67] Paranjape J N, Nair N G, Sikder S, et al. Adaptivesam: Towards efficient tuning of sam for surgical scene segmentation[C]. *Annual Conference on Medical Image Understanding and Analysis*, 2024: 187-201.
- [68] Hu E J, Shen Y, Wallis P, et al. Lora: Low-rank adaptation of large language models[J]. *arXiv preprint [arXiv:2106.09685](https://arxiv.org/abs/2106.09685)*, 2021.
- [69] Sung Y L, Cho J, Bansal M. Vi-adapter: Parameter-efficient transfer learning for vision-and-language tasks[C]. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022: 5227-5237.
- [70] Olabarriaga S D, Smeulders A W M. Interaction in the segmentation of medical images: A survey[J]. *Medical Image Analysis*, 2001, 5(2): 127-142.
- [71] Wong H E, Rakic M, Gutttag J, et al. Scribbleprompt: Fast and flexible interactive segmentation for any medical image[J]. *arXiv preprint [arXiv:2312.07381](https://arxiv.org/abs/2312.07381)*, 2023.
- [72] Wang G, Zuluaga M A, Pratt R, et al. Slic-Seg: A minimally interactive segmentation of the placenta from sparse and motion-corrupted fetal MRI in multiple views[J]. *Medical Image Analysis*, 2016, 34: 137-147.
- [73] Sakinis T, Milletari F, Roth H, et al. Interactive segmentation of medical images through fully convolutional neural networks[J]. *arXiv preprint [arXiv:1903.08205](https://arxiv.org/abs/1903.08205)*, 2019.
- [74] Du Y, Bai F, Huang T, et al. Segvol: Universal and interactive volumetric medical image segmentation[J]. *arXiv preprint [arXiv:2311.13385](https://arxiv.org/abs/2311.13385)*, 2023.
- [75] Lüddecke T, Ecker A. Image segmentation using text and image prompts[C]. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022: 7086-7096.
- [76] Grady L, Schiwietz T, Aharon S, et al. Random walks for interactive organ segmentation in two and three dimensions: Implementation and validation[C]. *Medical Image Computing and Computer-Assisted Intervention–MICCAI 2005, Springer*, 2005: 773-780.
- [77] Rajchl M, Lee M C H, Oktay O, et al. Deepcut: Object segmentation from bounding box annotations using convolutional neural networks[J]. *IEEE Transactions on Medical Imaging*, 2016, 36(2): 674-683.



## Appendix A: Demo and Code

Our code, model weights, and dataset are available at: <https://github.com/uni-medical/IMIS-Bench>.

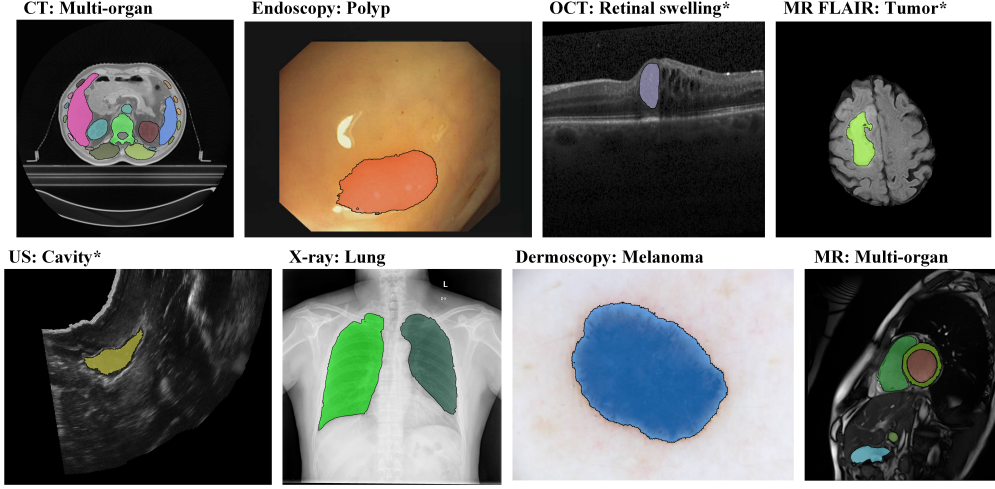


Figure 8: **Example predictions of IMIS-Net across different modalities and segmentation tasks.** "\*" Indicates that the corresponding image modality or segmentation task was not included in our training plan. Our model demonstrates its versatility by effectively handling multiple medical image modalities and performing various segmentation tasks, even on those that it has not previously encountered.

## Appendix B: IMed-361M Information and Availability

We have compiled 110 medical image segmentation datasets into a comprehensive, large-scale, multimodal, high-quality dataset named IMed-361M, making it openly accessible for interactive medical image segmentation. This dataset includes over 6.4 million images, 87.6 million ground truth annotations, and 273.4 million interactive masks (IMask), encompassing 14 image modalities and 204 segmentation targets. Fig.9 presents representative samples, while detailed category information is provided in Tab.3, which forms the basis for IMIS-Net’s text prompts. The modality and category details of these open-source and private datasets are displayed in Tab.4.

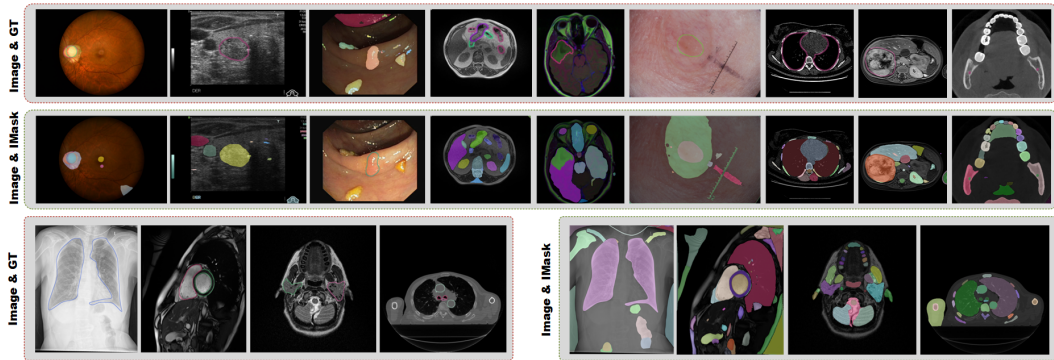


Figure 9: **IMed-361M:** A comprehensive dataset of multimodal medical images encompassing nearly all human organs and lesions, with interactive masks offering detailed, dense annotations.

**Division of Datasets.** Tab.2 summarizes the datasets used for training and testing our model. Each dataset was divided into 90% for training and 10% for testing, with 3D datasets split along the volume dimension. To ensure evaluation reliability, we limited the test set of each dataset to a maximum of

Table 3: IMed-361M dataset contains six anatomical categories: A (Abdomen), S (Skeleton), H (Head & Neck), T (Thorax), P (Pelvis), and L (Lesions). The symbol \* indicates that the target has left and right parts.

A01: Adrenal gland *	H07: Optic chiasm	S07: Scapula *	L01: Lung infections
A02: Aorta	H08: Pituitary gland	S08: Cervical spine (C1-C7)	L02: Liver tumor
A03: Autochthonous muscles*	H09: Brain stem	S09: Lumbar spine (L1-L6)	L03: Kidney tumor
A04: Colon	H10: Temporal lobe *	S10: Thoracic spine (T1-T13)	L04: Kidney cyst
A05: Duodenum	H11: Parotid gland *	T01: Esophagus	L05: Pleural effusion
A06: Gallbladder	H12: Ear (*, Inner, Middle)	T02: Atrium *	L06: Myocardial edema
A07: Iliac artery *	H13: Temporomandibular *	T03: Myocardium *	L07: Myocardial scars
A08: Iliac vein *	H14: Mandible *	T04: Ventricle *	L08: Necrosis
A09: Iliopsoas *	H15: Thyroid gland	T05: Lower lobe *	L09: Edema
A10: Inferior vena cava	H16: Submandibular gland *	T06: Middle lobe *	L10: Non enhancing tumor
A11: Kidney *	H17: Oral cavity	T07: Upper lobe *	L11: Enhancing tumor
A12: Liver	H18: Eustachian tube *	T08: Pulmonary artery	L12: Necrotic tumor core
A13: Pancreas	H19: Hippocampus *	T9: Trachea	L13: Peritumoral edema
A14: Portal and splenic veins	H20: Mastoid *	T10: Lung	L14: Myocardial infarction
A15: Small intestine	H21: Tympanic cavity *	T11: Heart	L15: No reflow
A16: Spleen	H22: Semicircular canal *	T12: Bronchus *	L16: Brain aneurysm
A17: Stomach	H23: Optic cup	T13: Breast *	L17: Neuroblastoma
A18: Spinal cord	H24: Optic disc	T14: Ascending aorta	L18: Prostate AFMS
A19: Rectum	H25: Larynx glottis	P01: Gluteus maximus *	L19: Hypoxic-ischemic
A20: Portal veins	H26: Larynx	P02: Gluteus medius *	L20: Breast tumor
A21: Large bowel	H27: Pharyngeal constrictor	P03: Gluteus minimus *	L21: Glioma
H01: Brain	S01: Clavicle *	P04: Bladder	L22: Thyroid nodule
H02: Face	S02: Femur *	P05: Prostate and uterus	L23: Skin lesion
H03: Airway	S03: Hip *	P06: Prostate	L24: Polyp
H04: Eye *	S04: Humerus *	P07: Testicle	P10: Prostatic urethra
H05: Crystalline lens *	S05: Rib (L&R, 1-12)	P08: Prostate peripheral zone	
H06: Optic nerve *	S06: Sacrum	P09: Prostate transition zone	

3,000 images for final model assessment. Additionally, we evaluated the model’s zero-shot capability using three external datasets: SegThor [129], TotalSegmentatorMRI [34], and ISLES<sup>9</sup>. Therefore, our training and test sets share the same data distribution, including modality and category.

Table 4: **Training and Test Datasets.** The following datasets were collected for training and validating IMIS-Net. Processed non-private datasets will be made publicly available for research purposes.

Dataset	Segmentation target	Modality	Category
SegRap2023 [130]	A18; H01,04-20,25-26; T01,09	CT	45
AbdomenAtlasMini1.0 [131]	A02,06,10-13,16-17	CT	9
AbdomenCT1K [39]	A11-13,16	CT	4
AMOS2022 [37]	A01-02,05-06,10-13,16-17; P04-05; T01	CT	15
BTCV [42]	A01-02,06,10-14,16-17; T01	CT	13
Colorectal_Liver_Metastases [132]	A12; L02	CT	2
Continuous_Registration_task1 [133]	T10	CT	1
COVID-19 CT scans [109, 110]	T10; L01	CT	1
CTSpine1K_Full [103]	S-08-10	CT	25
Finding-lungs-in-cTdata_3d [29]	T10	CT	1
FLARE21 [91]	A11-13,16	CT	4
FLARE22 [90]	A01-02,05-06,10-13,16-17; T01	CT	13
HCC-TACE-Seg <sup>10</sup>	A02,12,20; L02	CT	4
KiTS [66]	A11; L03	CT	2
KiTS2021 [67]	A11; L03-04	CT	3
KiTS2023 [134]	A11; L03-04	CT	3
Learn2Reg2022_AbdomenCTCT	A01-02,06,10-14,16-17; T01	CT	13
Learn2Reg2022_AbdomenMRCT	A11-12,16	CT	4
LITS [73]	A12; L02	CT	2
LUNA16 [72]	T10	CT	1

Continued on next page

<sup>9</sup><https://www.isles-challenge.org/>

<sup>10</sup><https://www.cancerimagingarchive.net/collection/hcc-tace-seg/>

Table 4 : Continued from previous page			
Dataset	Segmentation target	Modality	Category
MMWHS [76, 77, 78]	T02-04,08,14	CT	7
MSD_Liver [111]	A12; L02	CT	2
MSD_Spleen [111]	A16	CT	1
PleThora	T10; L05	CT	2
Prostate-AnatomicalEdge-Cases [111]	A19; S-02; P04,06	CT	5
SLIVER07 [127]	A12	CT	1
STACOM_SLAWT [13]	T02	CT	1
Totalsegmentator [34]	A01-17; H01-02; S-01-10; P01-04; T01-09	CT	104
VESSEL2012 <sup>11</sup>	T10	CT	1
Sz_cxr [18]	T10	X-Ray	1
WORD [28]	A01,04-06,11-13,15-17,19; S-02; P04; T01	CT	16
WenYi_CTA_Data	A02; T02,04	CTA	3
CMRxC-Motions	T03-04	MR-CMR	3
Myops2020 [114, 115]	T03-04; L06-07	MR	5
BraTS2013 [43, 44]	L08-11	MR-FLAIR	4
BraTS2015 [43, 44]	L08-11	MR-FLAIR	4
BraTS2018 [43, 45, 46]	L11-13	MR-FLAIR	3
BraTS2019 [43, 45, 46]	L11-13	MR-FLAIR	3
BraTS2020 [43, 45, 46]	L11-13	MR-FLAIR	3
BraTS2021 [45, 46, 47]	L11-13	MR-FLAIR	3
BraTS2023_GLI <sup>12</sup>	L08-09,11	MR-FLAIR	3
BraTS2023_MEN	L09-11	MR-FLAIR	3
BraTS2023_MET	L09-11	MR-FLAIR	3
BraTS2023_PED	L09-11	MR-FLAIR	3
BraTS2023_SSA	L09-11	MR-FLAIR	3
BraTS-TCGA-GBM <sup>13</sup>	L11-13	MR-FLAIR	3
BraTS-TCGA-LGG	L11-13	MR-FLAIR	3
SPPIN2023 <sup>14</sup>	L17	MR-T1GD	1
ATLAS2023	A12; L02	MR-T1W	2
CHAOS_Task_4 [81, 82, 83]	A11-12,16	MR-T1W	4
Learn2Reg2022_AbdomenMRCT	A11-12,16	MR-T1W	4
MSD_Prostate [111]	P08-09	MR-T2W	2
Myops2020 [114, 115]	T03-04,06-07	MR-T2W	5
CHAOS_Task_4 [81, 82, 83]	A11-12,16	MR-T2W	4
ISBI-MR-Prostate-2013 <sup>15</sup>	P06,08	MR-T2W	2
Prostate_MRI [111]	P06	MR-T2W	1
PROSTATEx-Seg-HiRes <sup>16</sup>	P06	MR-T2W	1
PROSTATEx-Seg-Zones	P08-10; L18	MR-T2W	4
u-RegPro <sup>17</sup>	P06	MR-T2W	1
ADAM2020	L16	MR-TOF	2
ACDC [1]	T03-04	MR	3
AMOS2022 [37]	A01-02,05-06,10-13,16-17; P04-05; T01	MR	15
EMIDEC [96]	T03-04; L14-15	MR	4
Heart_Seg_MRI [86]	T02	MR	1
MMWHS [76, 77, 78]	T02-04,08,14	MR	7
Mnms2	T03-04	MR	3
MSD_Heart [112]	T02	MR	1
PROMISE12 [118]	P06	MR	1
CETUS2014 <sup>18</sup>	T04	US	1
CuRIOUS2022_tumor	L19	US	1
TDSC-ABUS2023 <sup>19</sup>	L20	US	1
u-RegPro	P06	US	1
BraimMRI [128]	L21	MR	1
Brain-MRI	L21	MR	1
UW-Madison	A15,17,21	MR	3
DDTI	L22	US	1
drishti_gs_cup [97, 98]	H23	Fundus	1
drishti_gs_od [97, 98]	H24	Fundus	1
gamma [17, 23, 24]	H23-24	Fundus	2
ichallenge_adam_task2 []	H24	Fundus	1
PAPILA [135]	H23-24	Fundus	2
refuge2 [51, 52]	H23-24	Fundus	2
rimonedl	H23-24	Fundus	2
finding-lungs-in-cTdata_2d [29]	T10	CT	1

Continued on next page

<sup>11</sup><https://zenodo.org/records/8055066>

<sup>12</sup><https://www.synapse.org/Synapse:syn51156910/wiki/621282>

<sup>13</sup><https://www.cancerimagingarchive.net/analysis-result/brats-tcga-gbm/>

<sup>14</sup><https://github.com/myrthebuser/SPPIN2023>

<sup>15</sup><https://www.cancerimagingarchive.net/analysis-result/isbi-mr-prostate-2013/>

<sup>16</sup><https://www.cancerimagingarchive.net/analysis-result/prostatex-seg-hires/>

<sup>17</sup><https://muregpro.github.io/>

<sup>18</sup><https://www.creatis.insa-lyon.fr/Challenge/CETUS/>

<sup>19</sup><https://tdsc-abus2023.grand-challenge.org/>

Dataset	Segmentation target	Modality	Category
isic2016_task1 <sup>20</sup>	L23	Dermoscopy	1
isic2017_task1	L23	Dermoscopy	1
isic2018_task1	L23	Dermoscopy	1
ph2 [121]	L23	Dermoscopy	1
cvc_clinicedb [102]	L24	Endoscopy	1
endovis15 [95]	L24	Endoscopy	1
hyper-kvasir-segmented-images	L24	Endoscopy	1
kvasir_seg [8]	L24	Endoscopy	1
kvasir_seg_aliyun	L24	Endoscopy	1
kvasircapsule_seg [20]	L24	Endoscopy	1
sun_seg	L24	Endoscopy	1
SegRap2023 [130]	A18; H01,04-20,25-26; T01,09	CTA	45
Private1	A19; S-02; P04	CT	4
Private2	A01-17; H01-02; S-01-10; P01-04; T01-09	CT	104
Private3	A18; H09-10	CT	4
Private4	A01-17; H01; S-01-10; P01-04; T01-09	CT	104
Private5	H01	CT	1
Private6	H11, 13-14	CT	10
Private7	H15-17,25,27	CT	6
Private8	T01,09-11,18	CT	6
Private9	A11-12,16-17	CT	5
Private10	A01-17; H01-02; S-01-10; P01-04; T01-09	CT	104
Private11	A19; S-02; P04,06,07	CT	6
Private12	T13	CT	2
Private13	T02,04,09,12	CT	7
Private14	A01-19; H01-14; S-01-10; P01-04; T01-09	CT	130

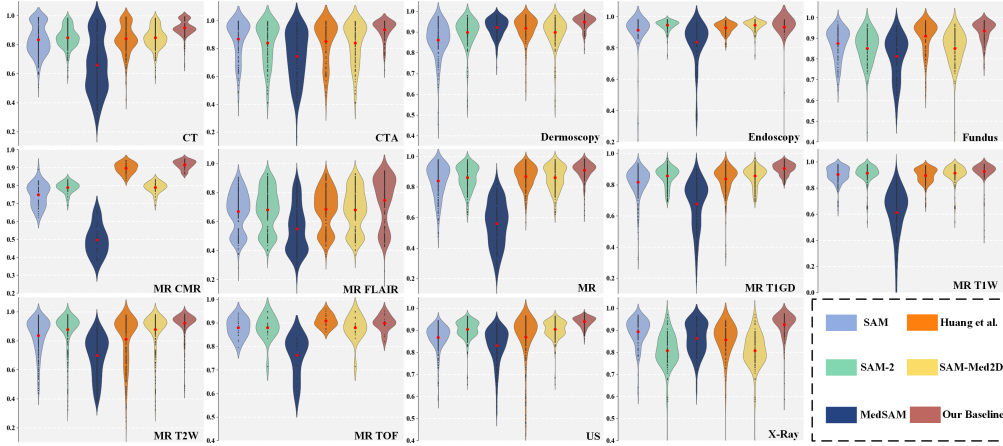


Figure 10: Comparison of segmentation performance of different methods in 14 medical image modalities, where the red points represent the means.

## Appendix C: Additional Experimental Analysis

**Comparative experiments with other models on different modalities.** Fig.10 shows the segmentation performance across 14 medical image modalities, with the red dots representing the average values. The methods proposed by MedSAM and Huang et al. utilize bounding box prompts, while the other models are evaluated based on better performance achieved from either bounding box prompts or three-click inputs. It can be observed that for medical modalities similar to natural images (e.g., dermoscopy and endoscopy), the performance of SAM and SAM-2 is comparable to the fine-tuned models, which validates the effectiveness of large-scale pretraining data. Additionally, our baseline model performs excellently across 12 modalities, with Dice scores exceeding 90%, and shows significant stability in modalities such as CTA, dermoscopy, ultrasound, and X-ray. These results suggest that directly using SAM or SAM-2 as a solution for medical image modalities with

<sup>20</sup><https://challenge.isic-archive.com/data/>

similar characteristics to natural images is feasible. However, for modalities that significantly differ from natural images, fine-tuning the base model can significantly improve segmentation performance.

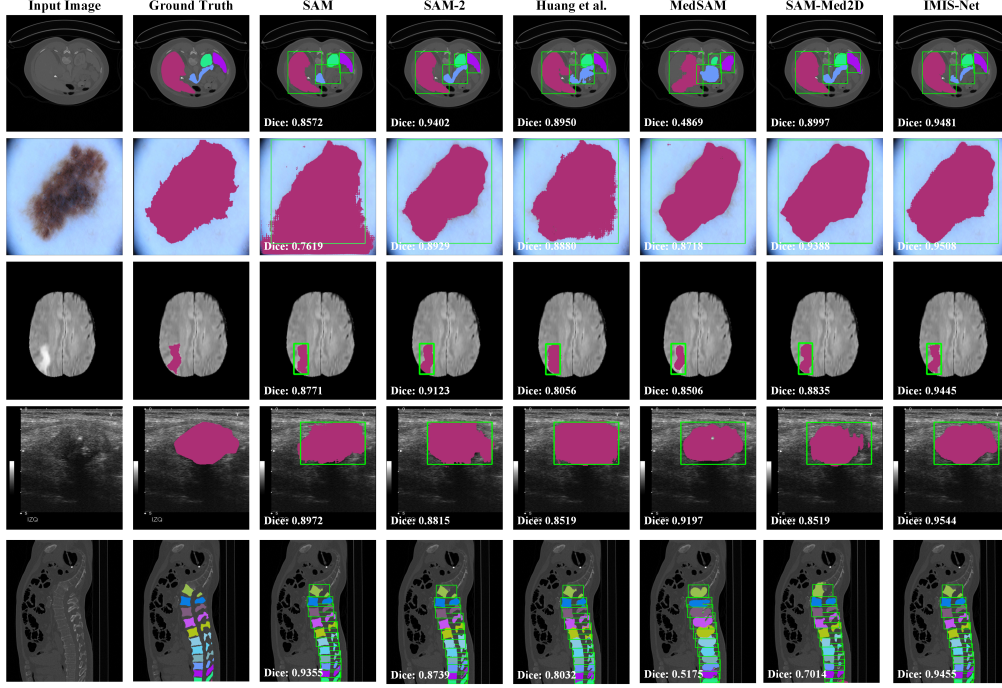


Figure 11: Simulate interactive segmentation results with identical bounding box coordinates for all model inputs.

**Visualization of interactive segmentation.** We assess interactive segmentation performance using Dice scores. The minimum enclosing bounding box of the ground truth serves as the model’s prompt input. Fig. 11 presents the prediction results generated by different models based on a single bounding box prompt. Due to the detailed spatial information provided by the bounding box, the Dice scores of various models are mostly above 0.8. In practical applications, this singular interactive approach may not directly meet user needs; hence, the IMIS method supports correcting predictions by providing additional click interactions. Our IMIS-Net still achieves high Dice segmentation. As shown in Fig. 12 and Fig. 13, we visualize the results of 5 simulated interactive experiments for SAM, SAM-2, and IMIS-Net.



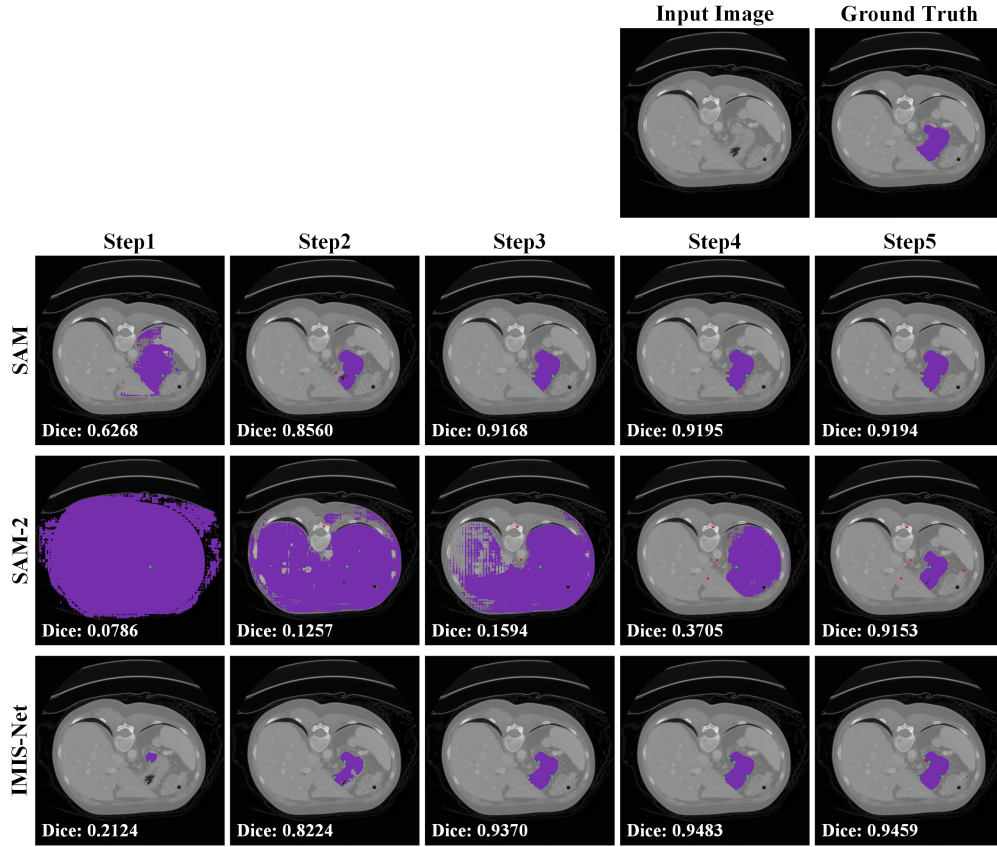


Figure 12: An interactive segmentation example of the stomach in CT images. SAM and SAM-2 typically require more prompts to achieve better results, while IMIS-Net achieves comparable performance with fewer interactions.

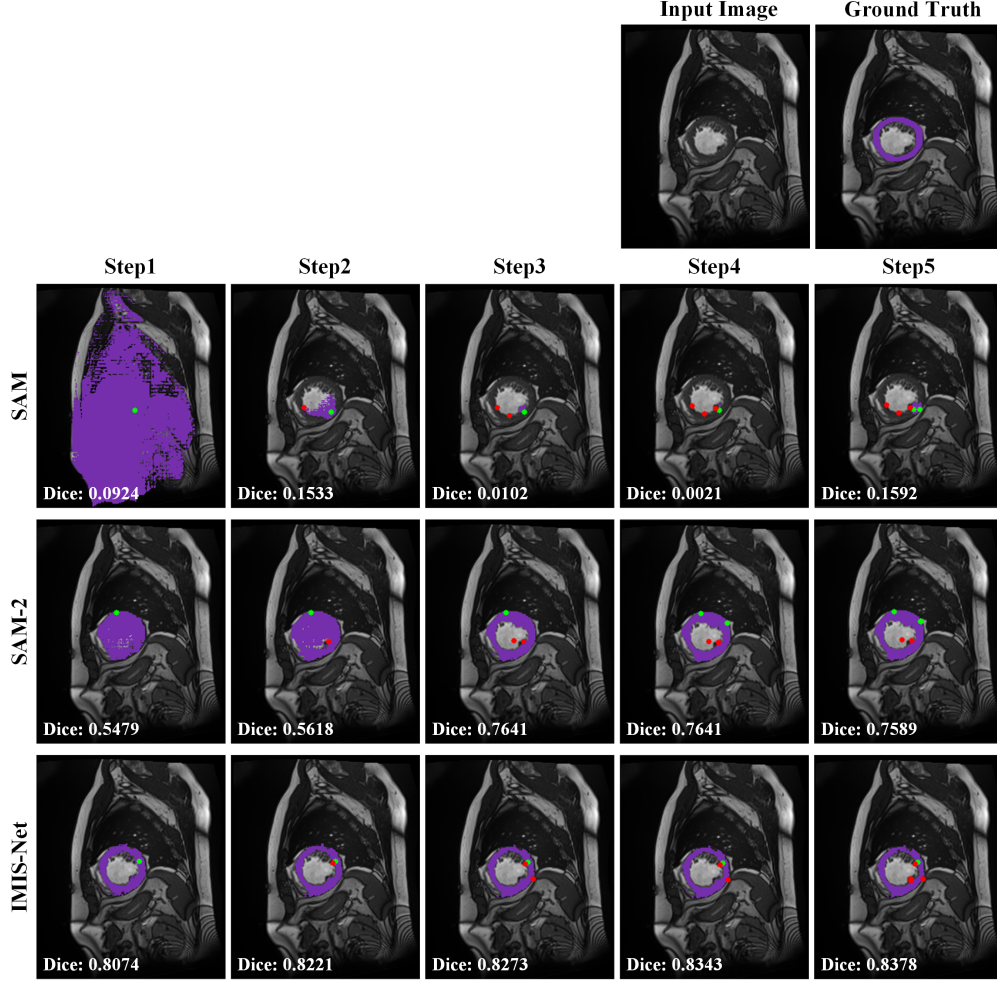


Figure 13: An interactive segmentation example of the cardiac myocardium in MR images. SAM performs poorly when dealing with annular myocardium, while SAM-2 and IMIS-Net are able to obtain predictions of the target area through multiple interactions. Our network consistently outperforms other methods.

## References

- [1] Bernard, Olivier, et al. "Deep learning techniques for automatic MRI cardiac multi-structures segmentation and diagnosis: is the problem solved?." *IEEE transactions on medical imaging* 37.11 (2018): 2514-2525.
- [2] Kaggle. Ultrasound Nerve Segmentation. <https://www.kaggle.com/competitions/ultrasound-nerve-segmentation/data>, 2016.
- [3] Saha, Anindo, Matin Hosseinzadeh, and Henkjan Huisman. "End-to-end prostate cancer detection in bpMRI via 3D CNNs: effects of attention mechanisms, clinical priori and decoupled false positive reduction." *Medical image analysis* 73 (2021): 102155.
- [4] Gutman, David, et al. "Skin lesion analysis toward melanoma detection: 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).
- [5] Codella, Noel CF, 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)." 2018 IEEE 15th international symposium on biomedical imaging (ISBI 2018). IEEE, 2018.
- [6] Tschandl, Philipp, Cliff Rosendahl, and Harald Kittler. "The HAM10000 dataset, a large collection of multi-source dermatoscopic images of common pigmented skin lesions." *Scientific data* 5.1 (2018): 1-9.
- [7] Codella, Noel, et al. "Skin lesion analysis toward melanoma detection 2018: A challenge hosted by the international skin imaging collaboration (isic)." *arXiv preprint arXiv:1902.03368* (2019).
- [8] Jha, Debesh, et al. "Kvasir-seg: A segmented polyp dataset." *MultiMedia modeling: 26th international conference, MMM 2020, Daejeon, South Korea, January 5–8, 2020, proceedings, part II* 26. Springer International Publishing, 2020.
- [9] Lu, Zhi, et al. "Evaluation of three algorithms for the segmentation of overlapping cervical cells." *IEEE journal of biomedical and health informatics* 21.2 (2016): 441-450.
- [10] Lu, Zhi, Gustavo Carneiro, and Andrew P. Bradley. "An improved joint optimization of multiple level set functions for the segmentation of overlapping cervical cells." *IEEE transactions on image processing* 24.4 (2015): 1261-1272.
- [11] Allan, Max, et al. "2017 robotic instrument segmentation challenge." *arXiv preprint arXiv:1902.06426* (2019).
- [12] Winzeck, Stefan, et al. "ISLES 2016 and 2017-benchmarking ischemic stroke lesion outcome prediction based on multispectral MRI." *Frontiers in neurology* 9 (2018): 679.
- [13] Karim, Rashed, et al. "Algorithms for left atrial wall segmentation and thickness–evaluation on an open-source CT and MRI image database." *Medical image analysis* 50 (2018): 36-53.
- [14] Commowick, Olivier, et al. "Multiple sclerosis lesions segmentation from multiple experts: The MICCAI 2016 challenge dataset." *Neuroimage* 244 (2021): 118589.
- [15] Martin Styner, Simon Warfield, Wiro Niessen, et al. MS Lesion Segmentation Challenge 2008. <http://www.ia.unc.edu/MSseg/index.html>, 2008.
- [16] Bharatha, Aditya, et al. "Evaluation of three-dimensional finite element-based deformable registration of pre-and intraoperative prostate imaging." *Medical physics* 28.12 (2001): 2551-2560.
- [17] Fu, Huazhu, et al. "Joint optic disc and cup segmentation based on multi-label deep network and polar transformation." *IEEE transactions on medical imaging* 37.7 (2018): 1597-1605.
- [18] Stirenko, Sergii, et al. "Chest X-ray analysis of tuberculosis by deep learning with segmentation and augmentation." 2018 IEEE 38th International Conference on Electronics and Nanotechnology (ELNANO). IEEE, 2018.
- [19] Shiraishi, Junji, et al. "Development of a digital image database for chest radiographs with and without a lung nodule: receiver operating characteristic analysis of radiologists' detection of pulmonary nodules." *American journal of roentgenology* 174.1 (2000): 71-74.
- [20] Jha, Debesh, et al. "Nanonet: Real-time polyp segmentation in video capsule endoscopy and colonoscopy." 2021 IEEE 34th International Symposium on Computer-Based Medical Systems (CBMS). IEEE, 2021.

- [21] Al-Dhabyani, Walid, et al. "Dataset of breast ultrasound images. Data Brief 28, 104863 (2020)." 2019,
- [22] Medimrg. RIM-ONE: Retinal Imaging - Medimrg. <http://medimrg.webs.ull.es/research/retinal-imaging/rim-one/>.
- [23] Orlando, José Ignacio, et al. "Refuge challenge: A unified framework for evaluating automated methods for glaucoma assessment from fundus photographs." *Medical image analysis* 59 (2020): 101570.
- [24] Fu, Huazhu, et al. "Age challenge: angle closure glaucoma evaluation in anterior segment optical coherence tomography." *Medical Image Analysis* 66 (2020): 101798.
- [25] Wang, Li, et al. "Benchmark on automatic six-month-old infant brain segmentation algorithms: the iSeg-2017 challenge." *IEEE transactions on medical imaging* 38.9 (2019): 2219-2230.
- [26] Sun, Yue, et al. "Multi-site infant brain segmentation algorithms: the iSeg-2019 challenge." *IEEE Transactions on Medical Imaging* 40.5 (2021): 1363-1376.
- [27] Gatidis, Sergios, et al. "A whole-body FDG-PET/CT dataset with manually annotated tumor lesions." *Scientific Data* 9.1 (2022): 601.
- [28] Luo, Xiangde, et al. "WORD: A large scale dataset, benchmark and clinical applicable study for abdominal organ segmentation from CT image." *Medical Image Analysis* 82 (2022): 102642.
- [29] Kuijf, Hugo J., et al. "Standardized assessment of automatic segmentation of white matter hyperintensities and results of the WMH segmentation challenge." *IEEE transactions on medical imaging* 38.11 (2019): 2556-2568.
- [30] Sekuboyina, Anjany, et al. "VerSe: a vertebrae labelling and segmentation benchmark for multi-detector CT images." *Medical image analysis* 73 (2021): 102166.
- [31] Löffler, Maximilian T., et al. "A vertebral segmentation dataset with fracture grading." *Radiology: Artificial Intelligence* 2.4 (2020): e190138.
- [32] Rudyanto, Rina D., et al. "Comparing algorithms for automated vessel segmentation in computed tomography scans of the lung: the VESSEL12 study." *Medical image analysis* 18.7 (2014): 1217-1232.
- [33] Baby, Mariya, and A. S. Jereesh. "Automatic nerve segmentation of ultrasound images." 2017 international conference of electronics, communication and aerospace technology (ICECA). Vol. 1. IEEE, 2017.
- [34] Wasserthal, Jakob, et al. "TotalSegmentator: robust segmentation of 104 anatomic structures in CT images." *Radiology: Artificial Intelligence* 5.5 (2023).
- [35] Kang, Qingbo, et al. "Thyroid nodule segmentation and classification in ultrasound images through intra-and inter-task consistent learning." *Medical image analysis* 79 (2022): 102443.
- [36] StructSeg2019 Grand Challenge Dataset. Retrieved from <https://structseg2019.grand-challenge.org/>
- [37] Ji, Yuanfeng, 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.
- [38] Zhang, Minghui, et al. "Multi-site, multi-domain airway tree modeling." *Medical image analysis* 90 (2023): 102957.
- [39] Ma, Jun, et al. "Abdomenct-1k: Is abdominal organ segmentation a solved problem?." *IEEE Transactions on Pattern Analysis and Machine Intelligence* 44.10 (2021): 6695-6714.
- [40] Xiong, Zhaohan, et al. "A global benchmark of algorithms for segmenting the left atrium from late gadolinium-enhanced cardiac magnetic resonance imaging." *Medical image analysis* 67 (2021): 101832.
- [41] Chen H, Zhao X, Dou, J. Carotid Vessel Wall Segmentation and Atherosclerosis Diagnosis Challenge. <https://vessel-wall-segmentation-2022.grand-challenge.org/>
- [42] Landman, Bennett, et al. "Miccai multi-atlas labeling beyond the cranial vault—workshop and challenge." *Proc. MICCAI Multi-Atlas Labeling Beyond Cranial Vault—Workshop Challenge*. Vol. 5. 2015.
- [43] Menze, Bjoern H., et al. "The multimodal brain tumor image segmentation benchmark (BRATS)." *IEEE transactions on medical imaging* 34.10 (2014): 1993-2024.

- [44] Kistler, Michael, et al. "The virtual skeleton database: an open access repository for biomedical research and collaboration." *Journal of medical Internet research* 15.11 (2013): e245.
- [45] Bakas, Spyridon, et al. "Advancing the cancer genome atlas glioma MRI collections with expert segmentation labels and radiomic features." *Scientific data* 4.1 (2017): 1-13.
- [46] Bakas, Spyridon, et al. "Identifying the best machine learning algorithms for brain tumor segmentation, progression assessment, and overall survival prediction in the BRATS challenge." *arXiv preprint arXiv:1811.02629* (2018).
- [47] Baid, Ujjwal, et al. "The rsna-asnr-miccai brats 2021 benchmark on brain tumor segmentation and radiogenomic classification." *arXiv preprint arXiv:2107.02314* (2021).
- [48] Avital, Itzik, et al. "Neural segmentation of seeding ROIs (sROIs) for pre-surgical brain tractography." *IEEE transactions on medical imaging* 39.5 (2019): 1655-1667.
- [49] Nelkenbaum, Ilya, et al. "Automatic segmentation of white matter tracts using multiple brain MRI sequences." 2020 IEEE 17th International Symposium on Biomedical Imaging (ISBI). IEEE, 2020.
- [50] CAD-PE Challenge | Developing a reference baseline for CAD-PE development. <http://www.cad-pe.org/>
- [51] Orlando, José Ignacio, et al. "Refuge challenge: A unified framework for evaluating automated methods for glaucoma assessment from fundus photographs." *Medical image analysis* 59 (2020): 101570.
- [52] Li, Fei, et al. "Development and clinical deployment of a smartphone-based visual field deep learning system for glaucoma detection." *NPJ digital medicine* 3.1 (2020): 123.
- [53] Xie, Ruitao, et al. "Pathological myopic image analysis with transfer learning." (2019).
- [54] Maier, Oskar, et al. "ISLES 2015-A public evaluation benchmark for ischemic stroke lesion segmentation from multispectral MRI." *Medical image analysis* 35 (2017): 250-269.
- [55] Cereda, Carlo W., et al. "A benchmarking tool to evaluate computer tomography perfusion infarct core predictions against a DWI standard." *Journal of Cerebral Blood Flow & Metabolism* 36.10 (2016): 1780-1789.
- [56] Hakim, Arsany, et al. "Predicting infarct core from computed tomography perfusion in acute ischemia with machine learning: Lessons from the ISLES challenge." *Stroke* 52.7 (2021): 2328-2337.
- [57] Hernandez Petzsche, Moritz R., et al. "ISLES 2022: A multi-center magnetic resonance imaging stroke lesion segmentation dataset." *Scientific data* 9.1 (2022): 762.
- [58] Li, Xiangyu, et al. "The state-of-the-art 3D anisotropic intracranial hemorrhage segmentation on non-contrast head CT: The INSTANCE challenge." *arXiv preprint arXiv:2301.03281* (2023).
- [59] Sirinukunwattana, Korsuk, David RJ Snead, and Nasir M. Rajpoot. "A stochastic polygons model for glandular structures in colon histology images." *IEEE transactions on medical imaging* 34.11 (2015): 2366-2378.
- [60] Sirinukunwattana, Korsuk, et al. "Gland segmentation in colon histology images: The glas challenge contest." *Medical image analysis* 35 (2017): 489-502.
- [61] Li, Xiangyu, et al. "Hematoma expansion context guided intracranial hemorrhage segmentation and uncertainty estimation." *IEEE Journal of Biomedical and Health Informatics* 26.3 (2021): 1140-1151.
- [62] He, Yuting, et al. "Meta grayscale adaptive network for 3D integrated renal structures segmentation." *Medical image analysis* 71 (2021): 102055.
- [63] He, Yuting, et al. "Dense biased networks with deep priori anatomy and hard region adaptation: Semi-supervised learning for fine renal artery segmentation." *Medical image analysis* 63 (2020): 101722.
- [64] Shao, Pengfei, et al. "Laparoscopic partial nephrectomy with segmental renal artery clamping: technique and clinical outcomes." *European urology* 59.5 (2011): 849-855.
- [65] Shao, Pengfei, et al. "Precise segmental renal artery clamping under the guidance of dual-source computed tomography angiography during laparoscopic partial nephrectomy." *European urology* 62.6 (2012): 1001-1008.



- [66] Heller, Nicholas, et al. "The kits19 challenge data: 300 kidney tumor cases with clinical context, ct semantic segmentations, and surgical outcomes." arXiv preprint [arXiv:1904.00445](https://arxiv.org/abs/1904.00445) (2019).
- [67] Zhao, Zhongchen, Huai Chen, and Lisheng Wang. "A coarse-to-fine framework for the 2021 kidney and kidney tumor segmentation challenge." International Challenge on Kidney and Kidney Tumor Segmentation. Cham: Springer International Publishing, 2021. 53-58.
- [68] Li, Lei, et al. "AtrialJSQnet: a new framework for joint segmentation and quantification of left atrium and scars incorporating spatial and shape information." Medical image analysis 76 (2022): 102303.
- [69] Li, Lei, et al. "Medical image analysis on left atrial LGE MRI for atrial fibrillation studies: A review." Medical image analysis 77 (2022): 102360.
- [70] Li, Lei, et al. "AtrialGeneral: domain generalization for left atrial segmentation of multi-center LGE MRIs." Medical Image Computing and Computer Assisted Intervention–MICCAI 2021: 24th International Conference, Strasbourg, France, September 27–October 1, 2021, Proceedings, Part VI 24. Springer International Publishing, 2021.
- [71] Pedrosa, João, et al. "LNDb: a lung nodule database on computed tomography." arXiv preprint [arXiv:1911.08434](https://arxiv.org/abs/1911.08434) (2019).
- [72] Setio, Arnaud Arindra Adiyoso, et al. "Validation, comparison, and combination of algorithms for automatic detection of pulmonary nodules in computed tomography images: the LUNA16 challenge." Medical image analysis 42 (2017): 1-13.
- [73] Bilic, Patrick, et al. "The liver tumor segmentation benchmark (lits)." Medical Image Analysis 84 (2023): 102680.
- [74] Carass, Aaron, et al. "Longitudinal multiple sclerosis lesion segmentation: resource and challenge." NeuroImage 148 (2017): 77-102.
- [75] Campello, Victor M., et al. "Multi-centre, multi-vendor and multi-disease cardiac segmentation: the M&Ms challenge." IEEE Transactions on Medical Imaging 40.12 (2021): 3543-3554.
- [76] Zhuang, Xiahai. "Multivariate mixture model for myocardial segmentation combining multi-source images." IEEE transactions on pattern analysis and machine intelligence 41.12 (2018): 2933-2946.
- [77] Zhuang, Xiahai, and Juan Shen. "Multi-scale patch and multi-modality atlases for whole heart segmentation of MRI." Medical image analysis 31 (2016): 77-87.
- [78] Luo, Xinzhe, and Xiahai Zhuang. " $\mathcal{X}$ -Metric: An N-Dimensional Information-Theoretic Framework for Groupwise Registration and Deep Combined Computing." IEEE Transactions on Pattern Analysis and Machine Intelligence 45.7 (2022): 9206-9224.
- [79] Mendrik, Adriënné M., et al. "MRBrainS challenge: online evaluation framework for brain image segmentation in 3T MRI scans." Computational intelligence and neuroscience 2015.1 (2015): 813696.
- [80] Van Ginneken, Bram, Tobias Heimann, and Martin Styner. "3D segmentation in the clinic: A grand challenge." MICCAI workshop on 3D segmentation in the clinic: a grand challenge. Vol. 1. 2007.
- [81] Kavur, A. Emre, et al. "CHAOS challenge-combined (CT-MR) healthy abdominal organ segmentation." Medical Image Analysis 69 (2021): 101950.
- [82] Kavur, Ali Emre, et al. "CHAOS-combined (CT-MR) healthy abdominal organ segmentation challenge data." Med. Image Anal 69 (2019): 101950.
- [83] Kavur, A. Emre, et al. "Comparison of semi-automatic and deep learning-based automatic methods for liver segmentation in living liver transplant donors." Diagnostic and Interventional Radiology 26.1 (2019): 11.
- [84] Wang, Shuo, et al. "The extreme cardiac MRI analysis challenge under respiratory motion (CMRxMotion)." arXiv preprint [arXiv:2210.06385](https://arxiv.org/abs/2210.06385) (2022).
- [85] Fang, Huihui, et al. "Adam challenge: Detecting age-related macular degeneration from fundus images." IEEE transactions on medical imaging 41.10 (2022): 2828-2847.
- [86] Tobon-Gomez, Catalina, et al. "Benchmark for algorithms segmenting the left atrium from 3D CT and MRI datasets." IEEE transactions on medical imaging 34.7 (2015): 1460-1473.

- [87] Pace, Danielle F., et al. "Interactive whole-heart segmentation in congenital heart disease." Medical Image Computing and Computer-Assisted Intervention–MICCAI 2015: 18th International Conference, Munich, Germany, October 5-9, 2015, Proceedings, Part III 18. Springer International Publishing, 2015.
- [88] Wang, Chuanbo, et al. "FUSeg: The foot ulcer segmentation challenge." Information 15.3 (2024): 140.
- [89] Payette, Kelly, et al. "An automatic multi-tissue human fetal brain segmentation benchmark using the fetal tissue annotation dataset." Scientific data 8.1 (2021): 167.
- [90] Ma, Jun, et al. "Unleashing the strengths of unlabeled data in pan-cancer abdominal organ quantification: the flare22 challenge." arXiv preprint [arXiv:2308.05862](https://arxiv.org/abs/2308.05862) (2023).
- [91] Ma, Jun, et al. "Fast and low-GPU-memory abdomen CT organ segmentation: the flare challenge." Medical Image Analysis 82 (2022): 102616.
- [92] Morozov, Sergey P., et al. "Mosmeddata: Chest ct scans with covid-19 related findings dataset." arXiv preprint [arXiv:2005.06465](https://arxiv.org/abs/2005.06465) (2020).
- [93] Aliyun Tianchi. Chest Image Dataset for Pneumothorax Segmentation. <https://tianchi.aliyun.com/dataset/83075>, 2020.
- [94] Ali, Sharib, et al. "Endoscopy disease detection challenge 2020." arXiv preprint [arXiv:2003.03376](https://arxiv.org/abs/2003.03376) (2020).
- [95] Bernal, Jorge, et al. "Comparative validation of polyp detection methods in video colonoscopy: results from the MICCAI 2015 endoscopic vision challenge." IEEE transactions on medical imaging 36.6 (2017): 1231-1249.
- [96] Lalonde, Alain, et al. "Deep learning methods for automatic evaluation of delayed enhancement-MRI. The results of the EMIDEC challenge." Medical Image Analysis 79 (2022): 102428.
- [97] Sivaswamy, Jayanthi, et al. "A comprehensive retinal image dataset for the assessment of glaucoma from the optic nerve head analysis." JSM Biomedical Imaging Data Papers 2.1 (2015): 1004.
- [98] Sivaswamy, Jayanthi, et al. "Drishti-gs: Retinal image dataset for optic nerve head (onh) segmentation." 2014 IEEE 11th international symposium on biomedical imaging (ISBI). IEEE, 2014.
- [99] Dorent, Reuben, et al. "CrossMoDA 2021 challenge: Benchmark of cross-modality domain adaptation techniques for vestibular schwannoma and cochlea segmentation." Medical Image Analysis 83 (2023): 102628.
- [100] Shusharina, Nadya, et al. "Cross-modality brain structures image segmentation for the radiotherapy target definition and plan optimization." Segmentation, Classification, and Registration of Multi-modality Medical Imaging Data: MICCAI 2020 Challenges, ABCs 2020, L2R 2020, TN-SCUI 2020, Held in Conjunction with MICCAI 2020, Lima, Peru, October 4–8, 2020, Proceedings 23. Springer International Publishing, 2021.
- [101] Hssayeni, Murtadha, et al. "Computed tomography images for intracranial hemorrhage detection and segmentation." Intracranial hemorrhage segmentation using a deep convolutional model. Data 5.1 (2020): 14.
- [102] Bernal, Jorge, et al. "WM-DOVA maps for accurate polyp highlighting in colonoscopy: Validation vs. saliency maps from physicians." Computerized medical imaging and graphics 43 (2015): 99-111.
- [103] Deng, Yang, et al. "CTSpine1K: a large-scale dataset for spinal vertebrae segmentation in computed tomography." arXiv preprint [arXiv:2105.14711](https://arxiv.org/abs/2105.14711) (2021).
- [104] Liu, Pengbo, et al. "Deep learning to segment pelvic bones: large-scale CT datasets and baseline models." International Journal of Computer Assisted Radiology and Surgery 16 (2021): 749-756.
- [105] Hogeweg, Laurens, et al. "Clavicle segmentation in chest radiographs." Medical image analysis 16.8 (2012): 1490-1502.
- [106] Cohen, Joseph Paul, et al. "Covid-19 image data collection: Prospective predictions are the future." arXiv preprint [arXiv:2006.11988](https://arxiv.org/abs/2006.11988) (2020).
- [107] Cohen, Joseph Paul, et al. "Covid-19 image data collection: Prospective predictions are the future." arXiv preprint [arXiv:2006.11988](https://arxiv.org/abs/2006.11988) (2020).
- [108] Roth, Holger R., et al. "Rapid artificial intelligence solutions in a pandemic—The COVID-19-20 Lung CT Lesion Segmentation Challenge." Medical image analysis 82 (2022): 102605.

- [109] Paiva, O. "Helping Radiologists To Help People In More Than 100 Countries!! Coronavirus Cases." CORONACASES. ORG, Ed., ed 11 (2020).
- [110] Ma, Jun, et al. "Toward data-efficient learning: A benchmark for COVID-19 CT lung and infection segmentation." *Medical physics* 48.3 (2021): 1197-1210.
- [111] Antonelli, Michela, et al. "The medical segmentation decathlon." *Nature communications* 13.1 (2022): 4128.
- [112] Simpson, Amber L., et al. "A large annotated medical image dataset for the development and evaluation of segmentation algorithms." *arXiv preprint [arXiv:1902.09063](https://arxiv.org/abs/1902.09063)* (2019).
- [113] Commowick, Olivier, et al. "Multiple sclerosis lesions segmentation from multiple experts: The MICCAI 2016 challenge dataset." *Neuroimage* 244 (2021): 118589.
- [114] Luo, Xinzhe, and Xiahai Zhuang. " $\mathcal{X}$ -Metric: An N-Dimensional Information-Theoretic Framework for Groupwise Registration and Deep Combined Computing." *IEEE Transactions on Pattern Analysis and Machine Intelligence* 45.7 (2022): 9206-9224.
- [115] Qiu, Junyi, et al. "MyoPS-Net: Myocardial pathology segmentation with flexible combination of multi-sequence CMR images." *Medical image analysis* 84 (2023): 102694.
- [116] Rister, Blaine, et al. "CT-ORG, a new dataset for multiple organ segmentation in computed tomography." *Scientific Data* 7.1 (2020): 381.
- [117] Saha, Anindo, Matin Hosseinzadeh, and Henkjan Huisman. "End-to-end prostate cancer detection in bpMRI via 3D CNNs: effects of attention mechanisms, clinical priori and decoupled false positive reduction." *Medical image analysis* 73 (2021): 102155.
- [118] Litjens, Geert, et al. "Evaluation of prostate segmentation algorithms for MRI: the PROMISE12 challenge." *Medical image analysis* 18.2 (2014): 359-373.
- [119] Pal, Jimut Bahan. "Holistic network for quantifying uncertainties in medical images." *International MICCAI Brainlesion Workshop*. Cham: Springer International Publishing, 2021.
- [120] Luo, Gongning, et al. "Efficient automatic segmentation for multi-level pulmonary arteries: The PARSE challenge." *arXiv preprint [arXiv:2304.03708](https://arxiv.org/abs/2304.03708)* (2023).
- [121] Mendonça, Teresa, et al. "PH 2-A dermoscopic image database for research and benchmarking." 2013 35th annual international conference of the IEEE engineering in medicine and biology society (EMBC). IEEE, 2013.
- [122] Zawacki, Anna, et al. "Siim-acr pneumothorax segmentation." (2019).
- [123] Liu, Quande, et al. "MS-Net: multi-site network for improving prostate segmentation with heterogeneous MRI data." *IEEE transactions on medical imaging* 39.9 (2020): 2713-2724.
- [124] Liu, Quande, Qi Dou, and Pheng-Ann Heng. "Shape-aware meta-learning for generalizing prostate MRI segmentation to unseen domains." *Medical Image Computing and Computer Assisted Intervention—MICCAI 2020: 23rd International Conference, Lima, Peru, October 4–8, 2020, Proceedings, Part II* 23. Springer International Publishing, 2020.
- [125] Jaeger, Stefan, et al. "Automatic tuberculosis screening using chest radiographs." *IEEE transactions on medical imaging* 33.2 (2013): 233-245.
- [126] Candemir, Sema, et al. "Lung segmentation in chest radiographs using anatomical atlases with nonrigid registration." *IEEE transactions on medical imaging* 33.2 (2013): 577-590.
- [127] Heimann, Tobias, et al. "Comparison and evaluation of methods for liver segmentation from CT datasets." *IEEE transactions on medical imaging* 28.8 (2009): 1251-1265.
- [128] Aliyun Tianchi. braimMRI Dataset. <https://tianchi.aliyun.com/dataset/127459>, 2022.
- [129] Lambert, Zoé, et al. "Segthor: Segmentation of thoracic organs at risk in ct images." 2020 Tenth International Conference on Image Processing Theory, Tools and Applications (IPTA). IEEE, 2020.
- [130] Luo, Xiangde, et al. "Segrap2023: A benchmark of organs-at-risk and gross tumor volume segmentation for radiotherapy planning of nasopharyngeal carcinoma." *arXiv preprint [arXiv:2312.09576](https://arxiv.org/abs/2312.09576)* (2023).

- [131] Qu, Chongyu, et al. "Abdomenatlas-8k: Annotating 8,000 CT volumes for multi-organ segmentation in three weeks." *Advances in Neural Information Processing Systems* 36 (2024).
- [132] Martin, Jack, et al. "Colorectal liver metastases: Current management and future perspectives." *World Journal of Clinical Oncology* 11.10 (2020): 761.
- [133] Marstal, Kasper, et al. "The continuous registration challenge: evaluation-as-a-service for medical image registration algorithms." *2019 IEEE 16th International Symposium on Biomedical Imaging (ISBI 2019)*. IEEE, 2019.
- [134] Myronenko, Andriy, et al. "Automated 3D Segmentation of Kidneys and Tumors in MICCAI KiTS 2023 Challenge." *International Challenge on Kidney and Kidney Tumor Segmentation*. Cham: Springer Nature Switzerland, 2023. 1-7.
- [135] Kovalyk, Oleksandr, et al. "PAPILA: Dataset with fundus images and clinical data of both eyes of the same patient for glaucoma assessment." *Scientific Data* 9.1 (2022): 291.