

Part-aware Prompted Segment Anything Model for Adaptive Segmentation

Chenhui Zhao

Department of Computer Science and Engineering, University of Michigan

chuizhao@umich.edu

Liyue Shen

Department of Electrical and Computer Engineering, University of Michigan

liyues@umich.edu

Reviewed on OpenReview: <https://openreview.net/forum?id=cCQKud5MFP¬eId=GEyhAYE7Q6>

Abstract

Precision medicine, such as patient-adaptive treatments assisted by medical image analysis, poses new challenges for segmentation algorithms in adapting to new patients, due to the large variability across different patients and the limited availability of annotated data for each patient. In this work, we propose a data-efficient segmentation algorithm, namely *Part-aware Prompted Segment Anything Model* (**P²SAM**). Without any model fine-tuning, P²SAM enables seamless adaptation to any new patients relying only on one-shot patient-specific data. We introduce a novel part-aware prompt mechanism to select multiple-point prompts based on the part-level features of the one-shot data, which can be extensively integrated into different promptable segmentation models, such as SAM and SAM 2. Moreover, to determine the optimal number of parts for each specific case, we propose a distribution-guided retrieval approach that further enhances the robustness of the part-aware prompt mechanism. P²SAM improves the performance by +8.0% and +2.0% mean Dice score for two different patient-adaptive segmentation applications, respectively. In addition, P²SAM also exhibits impressive generalizability in other adaptive segmentation tasks in the natural image domain, *e.g.*, +6.4% mIoU within personalized object segmentation task. The code is available at <https://github.com/Zch0414/p2sam>

1 Introduction

Advances in modern precision medicine and healthcare have emphasized the importance of patient-adaptive treatment (Hodson, 2016). For instance, in radiation therapy, the patient undergoing multi-fraction treatment would benefit from longitudinal medical data analysis that helps timely adjust treatment planning (Sonke et al., 2019). To facilitate the treatment procedure, such analysis demands timely and accurate automatic segmentation of tumors and critical organs from medical images, which has underscored the role of computer vision approaches for medical image segmentation tasks (Hugo et al., 2016; Jha et al., 2020). Despite the great progress made by previous works (Ronneberger et al., 2015; Isensee et al., 2021), their focus remains on improving the segmentation accuracy within a standard paradigm: trained on a large number of annotated data and evaluated on the *internal* validation set. However, patient-adaptive treatment presents unique challenges in adapting segmentation models to new patients: (1) the large variability across patients hinders direct model transfer, and (2) the limited availability of annotated training data for each patient prevents fine-tuning the model on a per-patient basis (Chen et al., 2023). Overcoming these obstacles requires a segmentation approach that can reliably adapt to *external* patients, in a data-efficient manner.

In this work, we address the unmet needs of the patient-adaptive segmentation by formulating it as an in-context segmentation problem, where the *context* is the prior data from a specific patient. Such data can be obtained in a standard clinical protocol (Chen et al., 2023), therefore will not burden clinician. To this end, we propose **P²SAM**: *Part-aware Prompted Segment Anything Model*. Leveraging the promptable segmentation mechanism inherent in Segment Anything Model (SAM) (Kirillov et al., 2023), our method seamlessly adapts to any *external* patients relying only on one-shot patient-specific prior data without re-

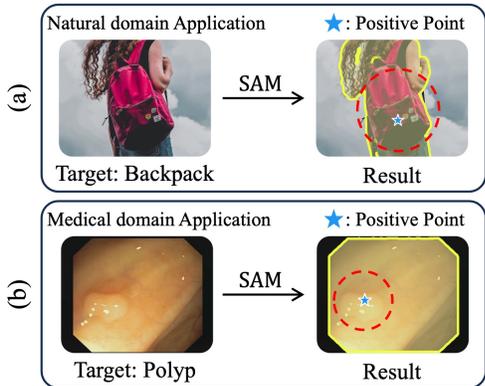


Figure 1: Illustration of SAM’s ambiguity property. The ground truth is circled by a red dashed circle; the predicted mask is depicted by a yellow solid line.

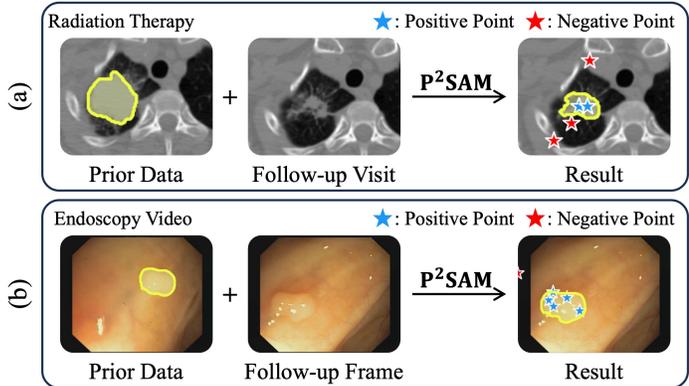


Figure 2: Illustration of two patient-adaptive segmentation tasks. P²SAM can segment the follow-up data by utilizing one-shot prior data as multiple-point prompts. Prior and predicted masks are depicted by a solid yellow line.

quiring additional training, thus in a data-efficient manner. Beyond patient-adaptive segmentation, P²SAM also demonstrates strong generalizability in other adaptive segmentation tasks in the natural image domain, such as personalized segmentation (Zhang et al., 2023) and one-shot segmentation (Liu et al., 2023).

In the original prompt mechanism of SAM (Kirillov et al., 2023), as illustrated in Figure 1, a single-point prompt may result in ambiguous prediction, indicating the limitation in both natural domain and medical domain applications (Zhang et al., 2023; Huang et al., 2024). To alleviate the ambiguity, following the statement in SAM, “*ambiguity is much rarer with multiple prompts*”, we propose a novel part-aware prompt mechanism that meticulously presents the prior data as multiple-point prompts based on part-level features. As illustrated in Figure 2, our method enables reliable adaptation to an *external* patient across various tasks with one-shot patient-specific prior data. To extract part-level features, we cluster the prior data into multiple parts in the feature space and computing the mean for each part. Then, we select multiple-point prompts based on the cosine similarity between these part-level features and the follow-up data. The proposed approach can be generalized to different promptable segmentation models that support the point modality, such as SAM and its successor, SAM 2 (Ravi et al., 2024). Here, we primarily utilize SAM as the backbone model, and SAM 2 will be integrated within the specific setting.

On the other hand, when the number of parts is set suboptimally, either more or less, the chance of encountering outlier prompts may increase. In the extreme, assigning all image patches to a single part produces an ambiguity-aware prompt (Zhang et al., 2023), whereas assigning each image patch to a different part yields many outlier prompts (Liu et al., 2023). Determining the optimal number of parts is non-trivial, as it may vary across different cases. Here, we introduced a novel distribution-guided retrieval approach to investigate the optimal number of parts required by each case. This retrieval approach is based on the distribution distance between the foreground feature of the prior image and the resulting feature obtained under the current part count. This principle is motivated by the fact that tumors and normal organs always lead to distinct feature distributions within medical imaging technologies (García-Figueiras et al., 2019).

With the aforementioned designs, P²SAM tackles a fundamental challenge—ambiguity—when adapting promptable segmentation models to specific applications. When ambiguity is not an issue, P²SAM enhances model generality by providing curated information. The key contributions of this work lie in three-fold:

1. We formulate the patient-adaptive segmentation as an in-context segmentation problem, resulting in a data-efficient segmentation approach, P²SAM, that requires only one-shot prior data and no model fine-tuning. P²SAM functions as a generic segmentation algorithm, enabling efficient and flexible adaptation across different domains, tasks, and models.
2. We propose a novel part-aware prompt mechanism that can select multiple-point prompts based on part-level features. Additionally, we introduce a distribution-guided retrieval approach to determine the optimal number of part-level features required by different cases. These designs significantly enhance the generalizability of promptable segmentation models.

3. Our method largely benefits real-world applications like patient-adaptive segmentation, one-shot segmentation, and personalized segmentation. Experiment results demonstrate that P²SAM improves the performance by +8.0% and +2.0% mean Dice score in two different patient-adaptive segmentation applications and achieves a new state-of-the-art result, *i.e.*, 95.7% mIoU on the personalized segmentation benchmark PerSeg.

2 Related Work

Segmentation Generalist. Over the past decade, various segmentation tasks including semantic segmentation (Strudel et al., 2021; Li et al., 2023a), instance segmentation (He et al., 2017; Li et al., 2022a), panoptic segmentation (Carion et al., 2020; Cheng et al., 2021; Li et al., 2022b), and referring segmentation (Li et al., 2023b; Zou et al., 2024) have been extensively explored for the image and video modalities. Motivated by the success of foundational language models (Radford et al., 2018; 2019; Brown et al., 2020; Touvron et al., 2023), the computer vision research community is increasingly paying attention to developing more generalized models that can tackle various vision or multi-modal tasks, or called foundation models (Li et al., 2022b; Oquab et al., 2023; Yan et al., 2023; Wang et al., 2023a;b; Kirillov et al., 2023). Notably, Segment Anything model (SAM) (Kirillov et al., 2023) and its successor, SAM 2 (Ravi et al., 2024) introduces a promptable model architecture, including the positive- and negative-point prompt; the box prompt; and the mask prompt. SAM and SAM 2 emerge with an impressive zero-shot interactive segmentation capability after pre-training on the large-scale dataset. The detail of SAM can be found in Appendix A.

Medical Segmentation. Given the remarkable generality of SAM and SAM 2, researchers within the medical image domain have been seeking to build foundational models for medical image segmentation (Wu et al., 2023; Wong et al., 2023; Wu & Xu, 2024; Zhang & Shen, 2024) in the same interactive fashion. To date, ScribblePrompt (Wong et al., 2023) and One-Prompt (Wu & Xu, 2024) introduce a new prompt modality—scribble—that provides a more flexible option for clinician usage. MedSAM (Ma et al., 2024a) fine-tunes SAM on an extensive medical dataset, demonstrating significant performance across various medical image segmentation tasks. Its successor (Ma et al., 2024b) incorporates SAM 2 to segment a 3D medical image volume as a video. However, these methods rely on clinician-provided prompts for promising segmentation performance. Moreover, whether these methods can achieve zero-shot performance as impressive as SAM and SAM 2 remains an open question that requires further investigation (Ma et al., 2024b).

In-Context Segmentation. The concept of in-context learning is first introduced as a new paradigm in natural language processing (Brown et al., 2020), allowing the model to adapt to unseen input patterns with a few prompts and examples, without the need to fine-tune the model. Similar ideas (Rakelly et al., 2018; Sonke et al., 2019; Li et al., 2023b) have been explored in segmentation tasks. For example, few-shot segmentation (Rakelly et al., 2018; Wang et al., 2019b; Liu et al., 2020; Leng et al., 2024) like PANet (Wang et al., 2019b), aims to segment new classes with only a few examples; in adaptive therapy (Sonke et al., 2019), several works (Elmahdy et al., 2020; Wang et al., 2020; Chen et al., 2023) attempt to adapt a segmentation model to new patients with limited patient-specific data, but these methods require model fine-tuning in different manners. Recent advancements, such as Painter (Wang et al., 2023a) and SegGPT (Wang et al., 2023b) pioneer novel in-context segmentation approaches, enabling the timely segmentation of images based on specified image-mask prompts. SEEM (Zou et al., 2024) further explores this concept by investigating different prompt modalities. More recently, PerSAM (Zhang et al., 2023) and Matcher (Liu et al., 2023) have utilized SAM to tackle few-shot segmentation through the in-context learning fashion. However, PerSAM prompts SAM with a single point prompt, causing ambiguity in segmentation results and therefore requires an additional fine-tuning strategy. Matcher samples multiple sets of point prompts but based on patch-level features. This mechanism makes Matcher dependent on DINOv2 (Oquab et al., 2023) to generate prompts, which is particularly pre-trained under a patch-level objective. Despite this, Matcher still generates a lot of outlier prompts, therefore relies on a complicated framework to filter the outlier results.

In this work, we address the patient-adaptive segmentation problem, also leveraging SAM’s promptable ability. Our prompt mechanism is based on part-level features, which will not cause ambiguity and are more robust than patch-level features. The optimal number of parts for each case is determined by a distribution-guided retrieval approach, further enhancing the generality of the part-aware prompt mechanism.

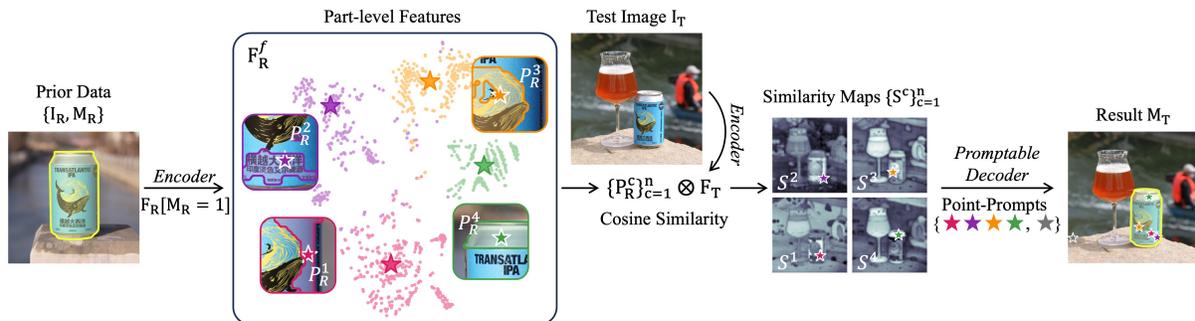


Figure 3: Illustration of the part-aware prompt mechanism. Masks are depicted by a yellow solid line. We first cluster foreground features in the reference image into part-level features. Then, we select multiple-point prompts based on the cosine similarity (\otimes in the figure) between these part-level features and target image features. A colorful star, matching the color of the corresponding part, denotes a positive-point prompt, while a gray star denotes a negative-point prompt. These prompts are subsequently fed into the promptable decoder to do prediction.

3 Method

In Section 3.1, we define the problem within the context of patient-adaptive segmentation. In Section 3.2, we present the proposed methodology, P²SAM, within a broader setting of adaptive segmentation. In Section 3.3, we introduce an optional fine-tuning strategy when adapting the backbone model to medical image domain is required.

3.1 Problem Setting

Our method aims to adapt a promptable segmentation model to *external* patients, with only one-shot patient-specific prior data. As shown in Figure 2, such data can be obtained in a standard clinical protocol, either from the initial visit of radiation therapy or the first frame of medical video. The prior data includes a reference image I_R and a mask M_R delineating the segmented object. Given a target image, I_T , our goal is to predict its mask M_T , without additional human annotation costs or model training burdens.

3.2 Methodology Overview

The setting described in Section 3.1 can be extended to other adaptive segmentation tasks in the natural image domain where the target image represents a new view or instance of the object depicted in the prior data. As shown in Figure 3, we illustrate our part-aware prompt mechanism using a natural image to clarify the significance of each part. Additional visualizations for parts in medical images are provided in Appendix D. Since no part-level definitions exist for the two diseases studied in this work, we refer these parts as data-driven parts.

Part-aware Prompt Mechanism. We utilize SAM (Kirillov et al., 2023) as the backbone model here, but our approach can be generalized to other promptable segmentation models that support the point prompt modality, such as SAM 2 (Ravi et al., 2024). Given the reference image-mask pair from the prior data, $\{I_R, M_R\}$, P²SAM first apply SAM’s *Encoder* to extract the visual features $F_R \in \mathbb{R}^{h \times w \times d}$ from the reference image I_R . Then, we utilize the reference mask M_R to select foreground features F_R^f ($F_R[M_R = 1]$) by:

$$F_R^f = \{F_{Rij} \mid M_{Rij} = 1, \forall (i, j) \in \mathcal{I}^{h \times w}\} \quad (1)$$

where $\mathcal{I}^{h \times w}$ is the spatial coordinate set of F_R . We cluster F_R^f with k -mean++ (Arthur et al., 2007) into n parts. Then, we obtain n part-level features $\{P_R^c\}_{c=1}^n \in \mathbb{R}^{n \times d}$ by computing the mean of each part. In Figure 3, we showcase an example of $n=4$. Each part-level feature P_R^c is represented by a colorful star in the foreground feature space. We further align the features of each part with pixels in the RGB space, thereby contouring the corresponding regions for each part in the image, respectively. We extract the

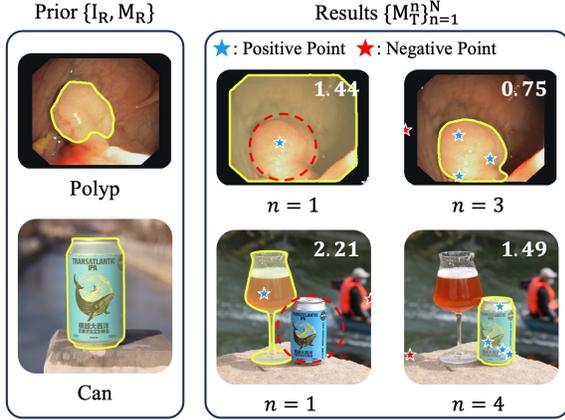


Figure 4: Illustration of P²SAM’s improvement. Wasserstein distances between the priors and results are shown in white.

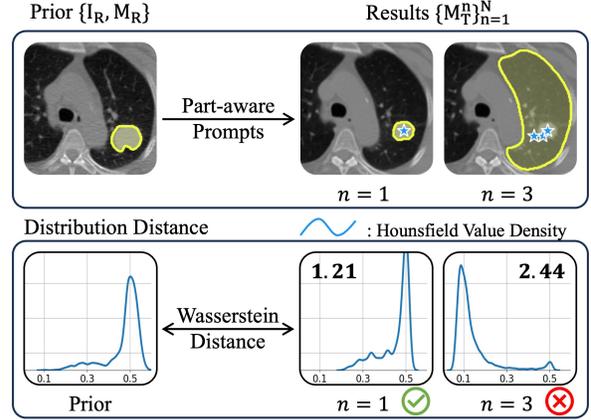


Figure 5: Illustration of the distribution-guided retrieval approach.

features $F_T \in \mathbb{R}^{h \times w \times d}$ from the target image I_T using the same *Encoder*, and compute similarity maps $\{S^c\}_{c=1}^n \in \mathbb{R}^{n \times h \times w}$ based on the cosine similarity between part-level features $\{P_R^c\}_{c=1}^n$ and F_T by:

$$S^c_{ij} = \frac{P_R^c \cdot F_{Tij}}{\|P_R^c\|_2 \cdot \|F_{Tij}\|_2} \quad (2)$$

We determine n positive-point prompts $\{Pos^c\}_{c=1}^n$ with the highest similarity score on each similarity map S^c . In Figure 3, each prompt Pos^c is depicted as a colorful star on the corresponding similarity map S^c .

For natural images, the background of the reference image and the target image may exhibit little correlation. Thus, following the approach in PerSAM (Zhang et al., 2023), we choose one negative-point prompt $\{Neg\}$ with the lowest score on the average similarity map $\frac{1}{n} \sum_{c=1}^n S^c$. $\{Neg\}$ is depicted as the gray star in Figure 3. However, for medical images, the background of the reference image is highly correlated with the background of the target image, usually both representing normal anatomical structures. As a result, in medical images, shown as Figure 2 in Section 1, we identify multiple negative-point prompts $\{Neg^c\}_{c=1}^n$ from the background. This procedure mirrors the selection of multiple positive-point prompts but we use background features F_R^b ($F_R[M_R = 0]$). Finally, we send both positive- and negative-point prompts into SAM’s *Promptable Decoder* and get the predicted mask M_T for the target image.

Distribution-Guided Retrieval Approach. Improvements of the part-aware prompt mechanism are illustrated in Figure 4. The proposed approach can naturally avoid the ambiguous prediction introduced by SAM (e.g., polyp) and also improve precision (e.g., can). However, this approach may occasionally result in outliers, as observed in the segmentation example in Figure 5, $n=3$. Therefore, we propose a distribution-guided retrieval approach to answer the question, “*How many part-level features should we choose for each case?*”. We assume the correct target foreground feature F_T^f ($F_T[M_T = 1]$), and the reference foreground feature F_R^f should belong to the same distribution. This assumption is grounded in the fact that tumors and normal organs will be reflected in distinct distributions by medical imaging technologies (García-Figueiras et al., 2019), also observed by the density of Hounsfield Unit value in Figure 5. To retrieve the optimal number of parts for a specific case, we first define N candidate part counts, and obtain N part-aware candidate segmentation results $\{M_T^n\}_{n=1}^N$. After that, we extract N sets of target foreground features $\{F_T^{f(n)}\}_{n=1}^N$. Following WGAN (Arjovsky et al., 2017), we utilize Wasserstein distance $\mathcal{D}_w(\cdot, \cdot)$ to measure the distribution distance between reference foreground features F_R^f and each set of target foreground features $F_T^{f(n)}$. We determine the optimal number of parts n by:

$$n = \arg \min_{n \in \{1, \dots, N\}} \mathcal{D}_w(F_R^f, F_T^{f(n)}), \quad (3)$$

where the details of $\mathcal{D}_w(\cdot, \cdot)$ can be found in Appendix F, Equation 5. The smaller distance value for the correct prediction in Figure 4 indicates this approach can be extended to multiple image modalities.

3.3 Adapt SAM to Medical Image Domain

Segment Anything Model (SAM) (Kirillov et al., 2023) is initially pre-trained on the SA-1B dataset. Despite the large scale, a notable domain gap persists between natural and medical images. In more realistic medical scenarios, clinic researchers could have access to certain public datasets (Aerts et al., 2015; Jha et al., 2020) tailored to specific applications, enabling them to fine-tune the model. Nevertheless, even after fine-tuning, the model can still be limited to generalize across various *external* patients from different institutions because of the large variability in patient population, demographics, imaging protocol, etc. P²SAM can then be flexibly plugged into the fine-tuned model to enhance robustness on unseen patients.

Specifically, when demanded, we utilize *internal* medical datasets (Aerts et al., 2015; Jha et al., 2020) to fine-tune SAM. We try full fine-tune, and Low-Rank adaptation (LoRA) (Hu et al., 2021) for further efficiency. During the fine-tuning, similar to Med-SA (Wu et al., 2023), we adhere closely to the interactive training strategy outlined in SAM to maintain the interactive ability. Details can be found in Appendix B. Then, we employ *external* datasets (Bernal et al., 2015; Hugo et al., 2016) obtained from various institutions to mimic new patient cases. Note that there is no further fine-tuning on these datasets.

4 Experiments

In Section 4.1, we introduce our experimental settings. In Section 4.2, we evaluate the quantitative results of our approach. In Section 4.3, we conducted several ablation studies to investigate our designs. In Section 4.4, we show qualitative results.

4.1 Experiment Settings

Dataset. We utilize a total of four medical datasets, including two *internal* datasets: The NSCLC-Radiomics dataset (Aerts et al., 2015), collected for non-small cell lung cancer (NSCLC) segmentation, contains data from 422 patients. Each patient has a computed tomography (CT) volume along with corresponding segmentation annotations. The Kvasir-SEG dataset (Jha et al., 2020), contains 1000 labeled endoscopy polyp images. Two *external* datasets from different institutions: The 4D-Lung dataset (Hugo et al., 2016), collected for longitudinal analysis, contains data from 20 patients, within which 13 patients underwent multiple visits, 3 to 8 visits for each patient. For each visit, a CT volume along with corresponding segmentation labels is available. The CVC-ClinicDB dataset (Bernal et al., 2015), contains 612 labeled polyp images selected from 29 endoscopy videos. During experiments, *internal* datasets serve as the training dataset to adapt SAM to the medical domain, while *external* datasets serve as unseen patient cases.

Patient-Adaptive Segmentation Tasks. We test P²SAM under two patient-adaptive segmentation tasks: NSCLC segmentation in the patient-adaptive radiation therapy and polyp segmentation in the endoscopy video. For NSCLC segmentation, medical image domain adaptation will be conducted on the *internal* dataset, NSCLC-Radiomics. For P²SAM, experiments are then carried out on the *external* dataset, 4D-Lung. We evaluate P²SAM on patients who underwent multiple visits during treatment. For each patient, we utilize the image-mask pair from the first visit as the patient-specific prior data. For polyp segmentation, domain adaptation will be conducted on *internal* dataset, Kvasir-SEG. For P²SAM, experiments are then carried out on *external* dataset, CVC-ClinicDB. For each video, we utilize the image-mask pair from the first stable frame as the patient-specific prior data.

Implementation Details. All experiments are conducted on A40 GPUs. For the NSCLC-Radiomics dataset, we extract 2-dimensional slices from the original computed tomography scans, resulting in a total of 7355 labeled images. As for the Kvasir-SEG dataset, we utilize all 1000 labeled images. We process two datasets following existing works (Hossain et al., 2019; Dumitru et al., 2023). Each dataset was randomly split into three subsets: training, validation, and testing, with an 80:10:10 percent ratio (patient-wise splitting for the NSCLC-Radiomics dataset to prevent data leak). The model is initialized with the SAM’s pre-trained weights and fine-tuned on the training splitting using the loss function proposed by SAM. We optimize the model by AdamW optimizer (Loshchilov & Hutter, 2017) ($\beta_1=0.9, \beta_2=0.999$), with a weight decay of 0.05. We further penalize the SAM’s encoder with a drop path of 0.1. We fine-tune the model for

Table 1: Results of NSCLC segmentation for patient-adaptive radiation therapy. We show the mean Dice score. *base*^{5.5M} indicates tuning 5.5M parameters of the base SAM on the NSCLC-Radiomics dataset before testing on the 4D-Lung dataset. † indicates training-free method; ‡ indicates the method using SAM.

Method	<i>Meta</i>	<i>LoRA</i>		<i>Full-Fine-Tune</i>	
	<i>huge</i> ^{0.0M}	<i>base</i> ^{5.5M}	<i>large</i> ^{5.9M}	<i>base</i> ^{93.8M}	<i>large</i> ^{312.5M}
<i>direct-transfer</i> [†]	-	56.10	57.83	58.18	61.11
<i>fine-tune</i>	-	52.11	32.55	55.27	53.85
PANet [†] (Wang et al., 2019b)	4.28	5.24	7.79	40.03	44.70
Matcher ^{†‡} (Liu et al., 2023)	13.28	50.81	50.88	59.52	57.67
PerSAM ^{†‡} (Zhang et al., 2023)	9.84	63.63	64.69	62.58	64.45
P ² SAM ^{†‡} (Ours)	28.52	64.38	67.00	66.68	67.23

Table 2: Results of polyp segmentation for endoscopy video. We show the mean Dice score for each method. *base*^{5.5M} indicates tuning 5.5M parameters of the base SAM on the Kvasir-SEG dataset before testing on the CVC-ClinicDB dataset. † indicates training-free method; ‡ indicates the method using SAM.

Method	<i>Meta</i>	<i>LoRA</i>		<i>Full-Fine-Tune</i>	
	<i>huge</i> ^{0.0M}	<i>base</i> ^{5.5M}	<i>large</i> ^{5.9M}	<i>base</i> ^{93.8M}	<i>large</i> ^{312.5M}
<i>direct-transfer</i> [†]	-	77.20	81.16	84.62	86.68
<i>fine-tune</i>	-	75.29	79.50	83.14	86.67
PANet [†] (Wang et al., 2019b)	38.22	44.61	55.48	75.99	86.48
Matcher ^{†‡} (Liu et al., 2023)	63.54	78.65	79.56	85.17	87.15
PerSAM ^{†‡} (Zhang et al., 2023)	45.82	79.02	81.63	85.74	87.88
P ² SAM ^{†‡} (Ours)	66.45	80.03	82.60	86.40	88.76

36 epochs on the NSCLC-Radiomics dataset and 100 epochs on the Kvasir-SEG dataset with a batch size of 4. The initial learning rate is $1e-4$, and the fine-tuning process is guided by cosine learning rate decay, with a linear learning rate warm-up over the first 10 percent epochs. More details are provided in Appendix C.

Summary. We test P²SAM on *external* datasets with three different SAM backbones: 1. SAM pre-trained on the SA-1B dataset (Kirillov et al., 2023), denoted as *Meta*. 2. SAM adapted on *internal* datasets with LoRA (Hu et al., 2021) and 3. full fine-tune, denoted as *LoRA* and *Full-Fine-Tune*, respectively. We compare P²SAM against various methods, including previous approaches such as the *direct-transfer*; *fine-tune* on the prior data (Wang et al., 2019a; Elmahdy et al., 2020; Wang et al., 2020; Chen et al., 2023); the one-shot segmentation method, PANet (Wang et al., 2019b); and concurrent methods that also utilize SAM, such as PerSAM (Zhang et al., 2023) and Matcher (Liu et al., 2023). For PANet, we utilize its align method for one-shot segmentation. For Matcher, we adopt its setting of FSS-1000 (Li et al., 2020). It is important to note that all baseline methods share the same backbone model as P²SAM does for fairness.

4.2 Quantitative Results

Patient-Adaptive Radiation Therapy. As shown in Table 1, on the 4D-Lung dataset (Hugo et al., 2016), P²SAM outperforms all other baselines across various backbones. Notably, when utilizing *Meta*, P²SAM can outperform Matcher by +15.24% and PerSAM by +18.68% mean Dice score. This highlights P²SAM’s superior adaptation to the out-of-domain medical applications. After domain adaptation, P²SAM outperforms the *direct-transfer* by +8.01%, Matcher by +11.60%, and PerSAM by +2.48% mean Dice score, demonstrating that P²SAM is a more effective method to enhance generalization on the *external* data.

Discussion. *fine-tune* is susceptible to overfitting with one-shot data, PANet fully depends on the encoder, and Matcher selects prompts based on patch-level features. These limitations prevent them from surpassing

Table 3: Comparison with existing baselines. \star indicates using a human-given box prompt during the inference time.

Method	4D-Lung	CVC-ClinicDB
<i>baseline</i>	69.00 \star	83.14
<i>direct-transfer</i>	61.11	86.68
P ² SAM	67.23	88.76

Table 5: Comparison with tracking methods. \star indicates utilizing *Full Fine-Tune*.

Method	4D-Lung	CVC-ClinicDB
AOT	-	62.34
P ² SAM	-	67.23
SAM 2	-	81.98
SAM 2 + P ² SAM	-	84.43
<i>label-propagation</i> \star	57.00	82.92
P ² SAM \star	67.23	88.76

Table 4: Results of one-shot semantic segmentation. We show the mean IoU score for each method. Note that all methods utilize SAM’s encoder for fairness.

Method	COCO-20 ⁱ	FSS-1000	LVIS-92 ⁱ	PerSeg
Matcher	25.1	82.1	12.6	90.2
PerSAM	23.0	71.2	11.5	89.3
P ² SAM (Ours)	26.0	82.4	13.7	95.7

Table 6: Ablation study for the number of parts n and the retrieval. Default settings are marked in Gray .

# parts (n)	CVC-ClinicDB		PerSeg	
	<i>w.o.</i>	<i>w. retrieval</i>	<i>w.o.</i>	<i>w. retrieval</i>
1 (PerSAM)	45.8	45.8	89.3	89.3
2	53.9	59.5	83.7	92.9
3	53.6	61.9	91.0	95.6
4	54.3	63.1	93.8	95.6
5	56.6	64.2	93.3	95.7

the *direct-transfer*. On the other hand, NSCLC segmentation remains a challenging task. We consider MedSAM (Ma et al., 2024a), which has been pre-trained on a large-scale medical image dataset, as a strong *baseline* method. In Table 3, MedSAM achieves a 69% mean dice score on the 4D-Lung dataset with a human-given box prompt at each visit, while P²SAM achieves comparable performance only with the ground truth provided at the first visit.

Endoscopy Video. As shown in Table 2, on the CVC-ClinicDB dataset (Bernal et al., 2015), P²SAM still achieves the best result across various backbones. When utilizing *Meta*, P²SAM can surpass Matcher by +2.91% and PerSAM by +20.63% mean Dice score. After domain adaptation, P²SAM can outperform *direct-transfer* by +2.03%, Matcher by +1.81% and PerSAM by +0.88% mean Dice score. Demonstrates P²SAM’s generality to various patient-adaptive segmentation tasks.

Discussion. All methods demonstrate improved performance in datasets like CVC-ClinicDB, which exhibit a smaller domain gap (Matsoukas et al., 2022) with SAM’s pre-training dataset. In Table 3, we compare our results with Sanderson & Matuszewski (2022), which is reported as the method achieving the best performance in Dumitru et al. (2023) under the same evaluation objective: trained on Kvasir-SEG dataset and tested on the CVC-ClinicDB dataset. Our *direct-transfer* has already surpassed this result, which can be attributed to the superior generality of SAM and our P²SAM can further improve the generalization.

On the other hand, we observe that P²SAM’s improvements over PerSAM become marginal after domain adaptation (*LoRA* and *Full Fine-Tune v.s. Meta*) on both datasets. This is because, as detailed in Appendix B, the ambiguity inherent in SAM, which is the primary limitation of PerSAM, is significantly reduced after fine-tuning on a dataset with a specific segmentation objective. Nevertheless, our method shows that providing multiple curated prompts can achieve further improvement.

Comparison with Tracking Algorithms. In Table 5, we additionally compared P²SAM with tracking algorithms: the *label-propagation* (Jabri et al., 2020), AOT (Yang et al., 2021), and SAM 2 (Ravi et al., 2024). On the 4D-Lung dataset, we only test algorithms with *Full Fine-Tune* due to the large domain gap (Matsoukas et al., 2022). P²SAM outperforms the *label-propagation*, as the discontinuity in sequential visits—where the interval between two CT scans can exceed a week—leads to significant changes in tumor position and features. On the CVC-ClinicDB dataset, dramatic content shifts within the narrow field of view can also lead to discontinuity. Despite this, SAM 2 achieves competitive results even without additional domain adaptation. However, as we have stated, P²SAM can be integrated into any promptable segmentation model. Indeed, we observe further improvements when applying P²SAM to SAM 2.

Table 7: Ablation study for the distribution distance measurement. Default settings are marked in Gray .

Algorithm	CVC-ClinicDB	PerSeg
<i>w.o.</i>	54.3	93.8
<i>Hungarian</i>	61.1	95.6
<i>Jensen-Shannon</i>	58.1	94.0
<i>Wasserstein</i>	63.1	95.6

Table 8: Ablation study for model sizes. \uparrow indicates the improvement when compared with the same size PerSAM. Default settings are marked in Gray .

Model	CVC-ClinicDB	PerSeg
PerSAM ^{huge}	45.8	89.3
P ² SAM ^{base}	55.1	90.0 _{26.0} \uparrow
P ² SAM ^{large}	63.8	95.6 _{9.0} \uparrow
P ² SAM ^{huge}	63.1	95.6 _{6.3} \uparrow

Existing One-shot Segmentation Benchmarks. To further demonstrate P²SAM can also be generalized to natural image domain, we evaluate its performance on existing one-shot semantic segmentation benchmarks: COCO-20ⁱ (Nguyen & Todorovic, 2019), FSS-1000 (Li et al., 2020), LVIS-92ⁱ (Liu et al., 2023), and a personalized segmentation benchmark, PerSeg (Zhang et al., 2023). We follow previous works (Zhang et al., 2023; Liu et al., 2023) for data pre-processing and evaluation. In Table 4, when utilizing SAM’s encoder, P²SAM outperforms concurrent works, Matcher and PerSAM, on all existing benchmarks. In addition, P²SAM can achieve a new state-of-the-art result, 95.7% mean IoU score, on the personalized segmentation benchmark PerSeg (Zhang et al., 2023).

4.3 Ablation Study

Ablation studies are conducted on the PerSeg dataset (Zhang et al., 2023) and CVC-ClinicDB dataset (Bernal et al., 2015) using *Meta*. We explore the effects of the number of parts in the part-aware prompt mechanism; the retrieval approach; distribution distance measurements in the retrieval approach; and the model size, which can be considered a proxy for representation capacity.

Number of Parts n . To validate the efficacy of the part-aware prompt mechanism, we establish a method without the retrieval approach. As shown in Table 6 (*w.o.* retrieval), for both datasets, even solely relying on the part-aware prompt mechanism, increasing the number of parts n enhances segmentation performance. When setting $n=5$, our part-aware prompt mechanism enhances performance by +10.7% mean Dice score on CVC-ClinicDB, +4.0% mean IoU score on PerSeg. These substantial improvements underscore the effectiveness of our part-aware prompt mechanism.

Retrieval Approach. The effectiveness of our retrieval approach is also shown in Table 6 (*w.* retrieval). When setting $n=5$, the retrieval approach enhances performance by +7.6% mean Dice score on the CVC-ClinicDB dataset, +2.4% mean IoU score on the PerSeg dataset. These substantial improvements show that our retrieval approach can retrieve an appropriate number of parts for different cases. Moreover, these suggest that we can initially define a wide range of part counts for retrieval, rather than tuning it meticulously as a hyperparameter.

Distribution Distance Measurements. The cornerstone of our retrieval approach lies in distribution distance measurements. To evaluate the efficacy of various algorithms, in Table 7, we juxtapose two distribution-related algorithms, namely *Wasserstein* distance (Rüschendorf, 1985) and *Jensen-Shannon* divergence (Menéndez et al., 1997), alongside a bipartite matching algorithm, *Hungarian* algorithm. Given foreground features from the reference image and the target image, we compute: 1. *Wasserstein* distance following the principles of WGAN (Arjovsky et al., 2017); 2. *Jensen-Shannon* divergence based on the first two principal components of each feature; 3. *Hungarian* algorithm after clustering these two sets of features into an equal number of parts. All algorithms exhibit improvements in segmentation performance compared to the *w.o.* retrieval baseline, while the *Wasserstein* distance is better in our context. Note that, the efficacy of the *Jensen-Shannon* divergence further corroborates our assumption that foreground features from the reference image and a correct target result should align in the same distribution, albeit it faces challenges when handling the high-dimensional data.

Model Size. In Table 8, we investigate the performance of different model sizes for our P²SAM, *i.e.*, *base*, *large*, and *huge*, which can alternatively be viewed as the representation capacity of different backbones. For

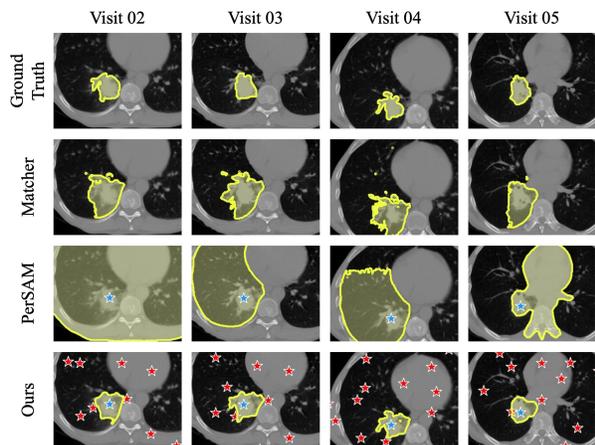


Figure 6: Qualitative results of NSCLC segmentation on the 4D-Lung dataset, with *Meta*.

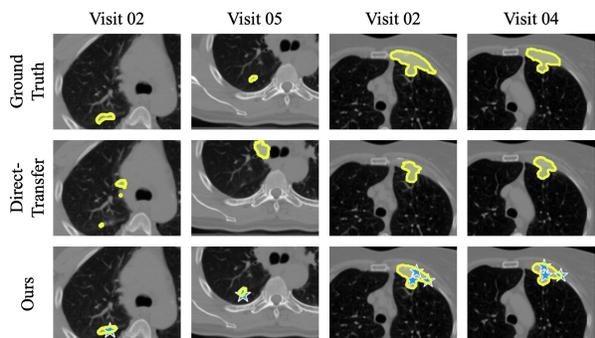


Figure 8: Qualitative results of NSCLC segmentation from two patients on the 4D-Lung dataset, with *Full-Fine-Tune*.

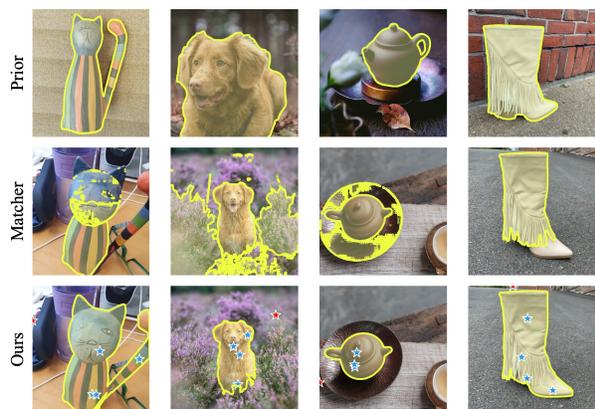


Figure 10: Qualitative results of personalized segmentation on the PerSeg dataset, compared with *Matcher*.

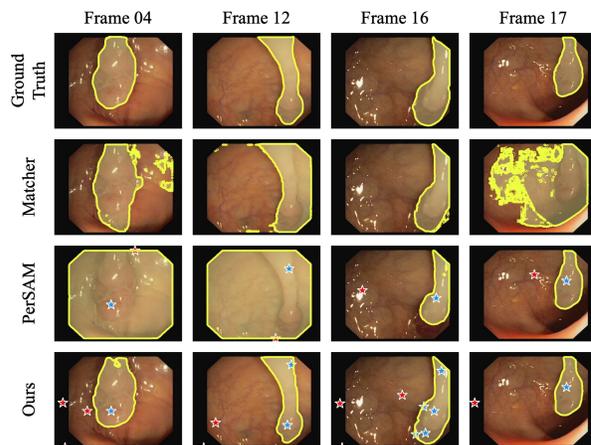


Figure 7: Qualitative results of polyp segmentation on the CVC-ClinicDB dataset, with *Meta*.

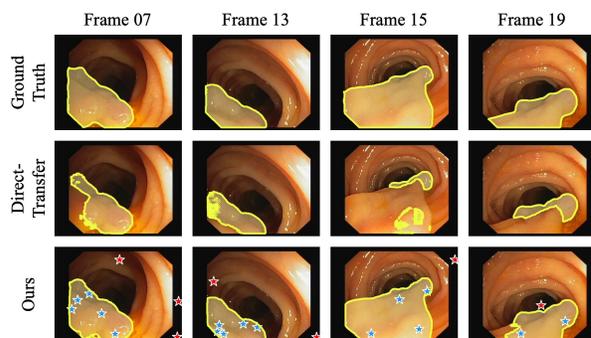


Figure 9: Qualitative results of polyp segmentation from one video on the CVC-ClinicDB dataset, with *Full-Fine-Tune*.

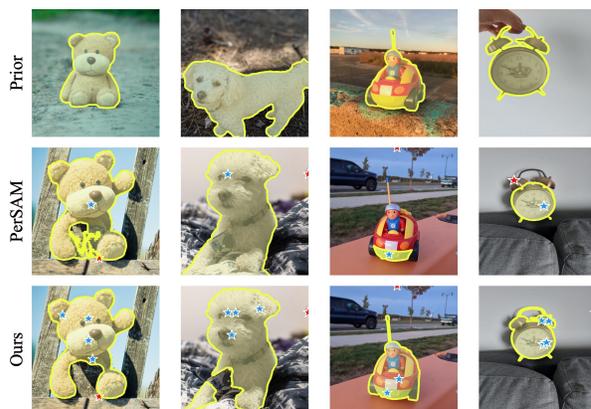


Figure 11: Qualitative results of personalized segmentation on the PerSeg dataset, compared with *PerSAM*.

the CVC-ClinicDB dataset, a larger model size does not necessarily lead to better results. This result aligns with current conclusions (Mazurowski et al., 2023; Huang et al., 2024): In medical image analysis, the *huge* SAM may occasionally be outperformed by the *large* SAM. On the other hand, for the PerSeg dataset, even utilizing the *base* SAM, P²SAM achieves higher accuracy compared to PerSAM with the *huge* SAM. These findings further underscore the robustness of P²SAM, particularly in scenarios where the model exhibits weaker representation, a circumstance more prevalent in medical image analysis.

4.4 Qualitative Results

Figure 6 and 7 showcase the advantage of P²SAM for out-of-domain applications. As shown in Figure 6, by presenting sufficient negative-point prompts, we enforce the model’s focus on the semantic target. Results in Figure 7 further summarize the benefits of our method: unambiguous segmentation and robust prompts selection. Our P²SAM can also improve the model’s generalization after domain adaptation. By providing precise foreground information, P²SAM enhances segmentation performance when the object is too small (*e.g.*, the first two columns in Figure 8) and when the segmentation is incomplete (*e.g.*, the last two columns in Figure 9). Figure 10 and 11 showcase the qualitative results on the PerSeg dataset, compared with Matcher and PerSAM respectively. The remarkable results demonstrate that P²SAM can generalize well to different domain applications.

5 Conclusion

We propose a data-efficient segmentation method, P²SAM, to solve the patient-adaptive segmentation problem. With a novel part-aware prompt mechanism and a distribution-guided retrieval approach, P²SAM can effectively integrate the patient-specific prior information into the current segmentation task. Beyond patient-adaptive segmentation, P²SAM demonstrates promising versatility in enhancing the backbone’s generalization across various levels: 1. At the domain level, P²SAM performs effectively in both medical and natural image domains. 2. At the task level, P²SAM enhances performance across different patient-adaptive segmentation tasks. 3. At the model level, P²SAM can be integrated into various promptable segmentation models, such as SAM, SAM 2, and custom fine-tuned SAM. In this work, to meet clinical requirements, we choose to adapt SAM to the medical imaging domain with public datasets. We opted not to adopt SAM 2, as it requires video data for fine-tuning, which is more costly. Additionally, treating certain patient-adaptive segmentation tasks as video tracking is inappropriate. In contrast, approaching patient-adaptive segmentation as an in-context segmentation problem offers a more flexible solution for various patient-adaptive segmentation tasks. Additional discussions can be found in appendix. We hope our work brings attention to the patient-adaptive segmentation problem within the research community.

Acknowledgments

The authors acknowledge support from University of Michigan MIDAS (Michigan Institute for Data Science) PODS Grant and University of Michigan MICDE (Michigan Institute for Computational Discovery and Engineering) Catalyst Grant, and the computing resource support from NSF ACCESS Program.

References

- HJWL Aerts, E Rios Velazquez, RT Leijenaar, Chintan Parmar, Patrick Grossmann, S Cavalho, Johan Bussink, René Monshouwer, Benjamin Haibe-Kains, Derek Rietveld, et al. Data from nsclc-radiomics. *The cancer imaging archive*, 2015.
- Michela Antonelli, Annika Reinke, Spyridon Bakas, Keyvan Farahani, Annette Kopp-Schneider, Bennett A Landman, Geert Litjens, Bjoern Menze, Olaf Ronneberger, Ronald M Summers, et al. The medical segmentation decathlon. *Nature communications*, 13(1):4128, 2022.
- Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein generative adversarial networks. In *International conference on machine learning*, pp. 214–223. PMLR, 2017.
- David Arthur, Sergei Vassilvitskii, et al. k-means++: The advantages of careful seeding. In *Soda*, volume 7, pp. 1027–1035, 2007.
- Jorge Bernal, F Javier Sánchez, Gloria Fernández-Esparrach, Debora Gil, Cristina Rodríguez, and Fernando Vilariño. Wm-dova maps for accurate polyp highlighting in colonoscopy: Validation vs. saliency maps from physicians. *Computerized medical imaging and graphics*, 43:99–111, 2015.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- Victor Ion Butoi, Jose Javier Gonzalez Ortiz, Tianyu Ma, Mert R Sabuncu, John Guttag, and Adrian V Dalca. Universeg: Universal medical image segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 21438–21451, 2023.
- Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *European conference on computer vision*, pp. 213–229. Springer, 2020.
- Yizheng Chen, Michael F Gensheimer, Hilary P Bagshaw, Santino Butler, Lequan Yu, Yuyin Zhou, Liyue Shen, Nataliya Kovalchuk, Murat Surucu, Daniel T Chang, et al. Patient-specific auto-segmentation on daily kvct images for adaptive radiotherapy. *International Journal of Radiation Oncology* Biology* Physics*, 2023.
- Bowen Cheng, Alex Schwing, and Alexander Kirillov. Per-pixel classification is not all you need for semantic segmentation. *Advances in Neural Information Processing Systems*, 34:17864–17875, 2021.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- Razvan-Gabriel Dumitru, Darius Peteleaza, and Catalin Craciun. Using duck-net for polyp image segmentation. *Scientific Reports*, 13(1):9803, 2023.
- Mohamed S Elmahdy, Tanuj Ahuja, Uulke A van der Heide, and Marius Staring. Patient-specific finetuning of deep learning models for adaptive radiotherapy in prostate ct. In *2020 IEEE 17th International Symposium on Biomedical Imaging (ISBI)*, pp. 577–580. IEEE, 2020.
- Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual object classes (voc) challenge. *International journal of computer vision*, 88:303–338, 2010.
- Roberto García-Figueiras, Sandra Baleato-González, Anwar R Padhani, Antonio Luna-Alcalá, Juan Antonio Vallejo-Casas, Evis Sala, Joan C Vilanova, Dow-Mu Koh, Michel Herranz-Carnero, and Herbert Alberto Vargas. How clinical imaging can assess cancer biology. *Insights into imaging*, 10:1–35, 2019.
- Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pp. 2961–2969, 2017.

- Richard Hodson. Precision medicine. *Nature*, 537(7619):S49–S49, 2016.
- Shahruk Hossain, Suhail Najeeb, Asif Shahriyar, Zaowad R. Abdullah, and M. Ariful Haque. A pipeline for lung tumor detection and segmentation from ct scans using dilated convolutional neural networks. In *ICASSP 2019 - 2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 1348–1352, 2019. doi: 10.1109/ICASSP.2019.8683802.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021.
- Yuhao Huang, Xin Yang, Lian Liu, Han Zhou, Ao Chang, Xinrui Zhou, Rusi Chen, Junxuan Yu, Jiongquan Chen, Chaoyu Chen, et al. Segment anything model for medical images? *Medical Image Analysis*, 92: 103061, 2024.
- Geoffrey D Hugo, Elisabeth Weiss, William C Sleeman, Salim Balik, Paul J Keall, Jun Lu, and Jeffrey F Williamson. Data from 4d lung imaging of nscL patients. *The Cancer Imaging Archive*, 10:K9, 2016.
- Fabian Isensee, Paul F Jaeger, Simon AA Kohl, Jens Petersen, and Klaus H Maier-Hein. nnu-net: a self-configuring method for deep learning-based biomedical image segmentation. *Nature methods*, 18(2): 203–211, 2021.
- Allan Jabri, Andrew Owens, and Alexei Efros. Space-time correspondence as a contrastive random walk. *Advances in neural information processing systems*, 33:19545–19560, 2020.
- Debesh Jha, Pia H Smedsrud, Michael A Riegler, Pål Halvorsen, Thomas de Lange, Dag Johansen, and Håvard D Johansen. Kvasir-seg: A segmented polyp dataset. In *MultiMedia Modeling: 26th International Conference, MMM 2020, Daejeon, South Korea, January 5–8, 2020, Proceedings, Part II 26*, pp. 451–462. Springer, 2020.
- Yuanfeng Ji, Haotian Bai, Chongjian Ge, Jie Yang, Ye Zhu, Ruimao Zhang, Zhen Li, Lingyan Zhang, Wanling Ma, Xiang Wan, et al. Amos: A large-scale abdominal multi-organ benchmark for versatile medical image segmentation. *Advances in neural information processing systems*, 35:36722–36732, 2022.
- Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. *arXiv preprint arXiv:2304.02643*, 2023.
- Tianang Leng, Yiming Zhang, Kun Han, and Xiaohui Xie. Self-sampling meta sam: Enhancing few-shot medical image segmentation with meta-learning. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pp. 7925–7935, 2024.
- Xiang Li, Tianhan Wei, Yau Pun Chen, Yu-Wing Tai, and Chi-Keung Tang. Fss-1000: A 1000-class dataset for few-shot segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 2869–2878, 2020.
- Xiang Li, Jinglu Wang, Xiao Li, and Yan Lu. Hybrid instance-aware temporal fusion for online video instance segmentation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pp. 1429–1437, 2022a.
- Xiang Li, Chung-Ching Lin, Yinpeng Chen, Zicheng Liu, Jinglu Wang, and Bhiksha Raj. Paintseg: Training-free segmentation via painting. *arXiv preprint arXiv:2305.19406*, 2023a.
- Xiang Li, Jinglu Wang, Xiaohao Xu, Xiao Li, Bhiksha Raj, and Yan Lu. Robust referring video object segmentation with cyclic structural consensus. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 22236–22245, 2023b.
- Yanhao Li, Hanzi Mao, Ross Girshick, and Kaiming He. Exploring plain vision transformer backbones for object detection. In *European Conference on Computer Vision*, pp. 280–296. Springer, 2022b.

- Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *Proceedings of the IEEE international conference on computer vision*, pp. 2980–2988, 2017.
- Yang Liu, Muzhi Zhu, Hengtao Li, Hao Chen, Xinlong Wang, and Chunhua Shen. Matcher: Segment anything with one shot using all-purpose feature matching. *arXiv preprint arXiv:2305.13310*, 2023.
- Yongfei Liu, Xiangyi Zhang, Songyang Zhang, and Xuming He. Part-aware prototype network for few-shot semantic segmentation. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part IX 16*, pp. 142–158. Springer, 2020.
- Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.
- Jun Ma, Yuting He, Feifei Li, Lin Han, Chenyu You, and Bo Wang. Segment anything in medical images. *Nature Communications*, 15(1):654, 2024a.
- Jun Ma, Sumin Kim, Feifei Li, Mohammed Baharoon, Reza Asakereh, Hongwei Lyu, and Bo Wang. Segment anything in medical images and videos: Benchmark and deployment. *arXiv preprint arXiv:2408.03322*, 2024b.
- Martin Maška, Vladimír Ulman, David Svoboda, Pavel Matula, Petr Matula, Cristina Ederra, Ainhoa Urbiola, Tomás España, Subramanian Venkatesan, Deepak MW Balak, et al. A benchmark for comparison of cell tracking algorithms. *Bioinformatics*, 30(11):1609–1617, 2014.
- Christos Matsoukas, Johan Fredin Haslum, Moein Sorkhei, Magnus Söderberg, and Kevin Smith. What makes transfer learning work for medical images: Feature reuse & other factors. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 9225–9234, 2022.
- Maciej A Mazurowski, Haoyu Dong, Hanxue Gu, Jichen Yang, Nicholas Konz, and Yixin Zhang. Segment anything model for medical image analysis: an experimental study. *Medical Image Analysis*, 89:102918, 2023.
- María Luisa Menéndez, JA Pardo, L Pardo, and MC Pardo. The jensen-shannon divergence. *Journal of the Franklin Institute*, 334(2):307–318, 1997.
- Fausto Milletari, Nassir Navab, and Seyed-Ahmad Ahmadi. V-net: Fully convolutional neural networks for volumetric medical image segmentation. In *2016 fourth international conference on 3D vision (3DV)*, pp. 565–571. Ieee, 2016.
- Khoi Nguyen and Sinisa Todorovic. Feature weighting and boosting for few-shot segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 622–631, 2019.
- Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning robust visual features without supervision. *arXiv preprint arXiv:2304.07193*, 2023.
- Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. Improving language understanding by generative pre-training. 2018.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- Kate Rakelly, Evan Shelhamer, Trevor Darrell, Alyosha Efros, and Sergey Levine. Conditional networks for few-shot semantic segmentation. 2018.
- Nikhila Ravi, Valentin Gabeur, Yuan-Ting Hu, Ronghang Hu, Chaitanya Ryali, Tengyu Ma, Haitham Khedr, Roman Rädle, Chloe Rolland, Laura Gustafson, et al. Sam 2: Segment anything in images and videos. *arXiv preprint arXiv:2408.00714*, 2024.

- Olaf Ronneberger, Philipp Fischer, and Thomas 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*, pp. 234–241. Springer, 2015.
- Ludger Rüschendorf. The wasserstein distance and approximation theorems. *Probability Theory and Related Fields*, 70(1):117–129, 1985.
- Edward Sanderson and Bogdan J Matuszewski. Fcn-transformer feature fusion for polyp segmentation. In *Annual conference on medical image understanding and analysis*, pp. 892–907. Springer, 2022.
- Jan-Jakob Sonke, Marianne Aznar, and Coen Rasch. Adaptive radiotherapy for anatomical changes. In *Seminars in radiation oncology*, volume 29, pp. 245–257. Elsevier, 2019.
- Robin Strudel, Ricardo Garcia, Ivan Laptev, and Cordelia Schmid. Segmenter: Transformer for semantic segmentation. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 7262–7272, 2021.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- A Vaswani. Attention is all you need. *Advances in Neural Information Processing Systems*, 2017.
- Chuang Wang, Neelam Tyagi, Andreas Rimner, Yu-Chi Hu, Harini Veeraraghavan, Guang Li, Margie Hunt, Gig Mageras, and Pengpeng Zhang. Segmenting lung tumors on longitudinal imaging studies via a patient-specific adaptive convolutional neural network. *Radiotherapy and Oncology*, 131:101–107, 2019a.
- Chuang Wang, Sadegh R Alam, Siyuan Zhang, Yu-Chi Hu, Saad Nadeem, Neelam Tyagi, Andreas Rimner, Wei Lu, Maria Thor, and Pengpeng Zhang. Predicting spatial esophageal changes in a multimodal longitudinal imaging study via a convolutional recurrent neural network. *Physics in Medicine & Biology*, 65(23):235027, 2020.
- Kaixin Wang, Jun Hao Liew, Yingtian Zou, Daquan Zhou, and Jiashi Feng. Panet: Few-shot image semantic segmentation with prototype alignment. In *proceedings of the IEEE/CVF international conference on computer vision*, pp. 9197–9206, 2019b.
- Xinlong Wang, Wen Wang, Yue Cao, Chunhua Shen, and Tiejun Huang. Images speak in images: A generalist painter for in-context visual learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 6830–6839, 2023a.
- Xinlong Wang, Xiaosong Zhang, Yue Cao, Wen Wang, Chunhua Shen, and Tiejun Huang. Seggpt: Segmenting everything in context. *arXiv preprint arXiv:2304.03284*, 2023b.
- Hallee E Wong, Marianne Rakic, John Guttag, and Adrian V Dalca. Scribbleprompt: Fast and flexible interactive segmentation for any medical image. *arXiv preprint arXiv:2312.07381*, 2023.
- Junde Wu and Min Xu. One-prompt to segment all medical images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 11302–11312, 2024.
- Junde Wu, Rao Fu, Huihui Fang, Yuanpei Liu, Zhaowei Wang, Yanwu Xu, Yueming Jin, and Tal Arbel. Medical sam adapter: Adapting segment anything model for medical image segmentation. *arXiv preprint arXiv:2304.12620*, 2023.
- Bin Yan, Yi Jiang, Jiannan Wu, Dong Wang, Ping Luo, Zehuan Yuan, and Huchuan Lu. Universal instance perception as object discovery and retrieval. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 15325–15336, 2023.
- Zongxin Yang, Yunchao Wei, and Yi Yang. Associating objects with transformers for video object segmentation. *Advances in Neural Information Processing Systems*, 34:2491–2502, 2021.

Anqi Zhang, Guangyu Gao, Jianbo Jiao, Chi Liu, and Yunchao Wei. Bridge the points: Graph-based few-shot segment anything semantically. *Advances in Neural Information Processing Systems*, 37:33232–33261, 2025.

Renrui Zhang, Zhengkai Jiang, Ziyu Guo, Shilin Yan, Junting Pan, Hao Dong, Peng Gao, and Hongsheng Li. Personalize segment anything model with one shot. *arXiv preprint arXiv:2305.03048*, 2023.

Yichi Zhang and Zhenrong Shen. Unleashing the potential of sam2 for biomedical images and videos: A survey. *arXiv preprint arXiv:2408.12889*, 2024.

Xueyan Zou, Jianwei Yang, Hao Zhang, Feng Li, Linjie Li, Jianfeng Wang, Lijuan Wang, Jianfeng Gao, and Yong Jae Lee. Segment everything everywhere all at once. *Advances in Neural Information Processing Systems*, 36, 2024.

Appendix

- A: SAM Overview
- B: SAM Adaptation Details
- C: Test Implementation Details
- D: Additional Visualizations
- E: Discussions
- F: Equations

A SAM Overview

Segment Anything Model (SAM) (Kirillov et al., 2023) comprises three main components: an image encoder, a prompt encoder, and a mask decoder, denoted as Enc_I , Enc_P , and Dec_M . As a promptable segmentation model, SAM takes an image I and a set of human-given prompts P as input. SAM predicts segmentation masks Ms by:

$$Ms = Dec_M(Enc_I(I), Enc_P(P)) \quad (4)$$

During training, SAM supervises the mask prediction with a linear combination of focal loss (Lin et al., 2017) and dice loss (Milletari et al., 2016) in a 20:1 ratio. When only a single prompt is provided, SAM generates multiple predicted masks. However, SAM backpropagates from the predicted mask with the lowest loss. Note that SAM returns only one predicted mask when presented with multiple prompts simultaneously.

Enc_I and Dec_M primarily employ the Transformer (Vaswani, 2017; Dosovitskiy et al., 2020) architecture. Here, we provide details on components in Enc_P . Enc_P supports three prompt modalities as input: the point, box, and mask logit. The positive- and negative-point prompts are represented by two learnable embeddings, denoted as E_{pos} and E_{neg} , respectively. The box prompt comprises two learnable embeddings representing the left-up and right-down corners of the box, denoted as E_{up} and E_{down} . In cases where neither the point nor box prompt is provided, another learnable embedding $E_{not-a-point}$ is utilized. If available, the mask prompt is encoded by a stack of convolution layers, denoted as E_{mask} ; otherwise, it is represented by a learnable embedding $E_{not-a-mask}$.

SAM employs an interactive training strategy. In the first iteration, either a positive-point prompt, represented by E_{pos} , or a box prompt, represented by $\{E_{up}, E_{down}\}$, is randomly selected with equal probability from the ground truth mask. Since there is no mask prompt in the first iteration, E_{pos} or $\{E_{up}, E_{down}\}$ is combined with $E_{not-a-mask}$ and fed into Dec_M . In the follow-up iterations, subsequent positive- and negative-point prompts are uniformly selected from the error region between the predicted mask and the ground truth mask. SAM additionally provides the mask logit prediction from the previous iteration as a supplement prompt. As a result, $\{E_{pos}, E_{neg}, E_{mask}\}$ is fed into Dec_M during each iteration. There are 11 total iterations: one sampled initial input prompt, 8 iteratively sampled points, and two iterations where only the mask prediction from the previous iteration is supplied to the model.

B SAM Adaptation Details

In Section 3.3, we propose to adapt SAM to the medical image domain when it is needed, with full fine-tune (*Full-Fine-Tune*) and LoRA (Hu et al., 2021) (*LoRA*). For *Full-Fine-Tune*, we fine-tune all parameters in SAM backbone. For *LoRA*, we insert the LoRA module in the image encoder Enc_I and only fine-tune parameters in the LoRA module and the mask decoder Dec_M . Our fine-tuning objectives are as follows:

1. The model can accurately predict a mask even if no prompt is provided.
2. The model can predict an exact mask even if only one prompt is given.

3. The model maintains promptable ability.

The training strategy outlined in SAM cannot satisfy all these three requirements: 1. The mask decoder Dec_M is not trained to handle scenarios where no prompt is given. 2. The approach to resolving the ambiguous prompt by generating multiple results is redundant as we have a well-defined segmentation objective. Despite that, we find a simple modification can meet all our needs:

1. In the initial iteration, we introduce a scenario where no prompt is provided to SAM. As a result, $\{E_{\text{not-a-point}}, E_{\text{not-a-mask}}\}$ is fed into Dec_M in the first iteration.
2. To prevent $E_{\text{not-a-point}}$ and $E_{\text{not-a-mask}}$ from introducing noise when human-given prompts are available, we stop their gradients in every iteration.
3. We ensure that SAM always returns an exact predicted mask. As a result, the ambiguity property does not exist in the model after fine-tuning.

C Test Implementation Details

Table 9: Retrieval range for the COCO-20ⁱ, FSS-1000, LVIS-92ⁱ, PerSeg dataset. **Blue** indicates the retrieval range for positive-point prompts. **Red** indicates the retrieval range for negative-point prompts.

COCO-20 ⁱ	FSS-1000	LVIS-92 ⁱ	PerSeg
1, 5-10 / 1	1-5 / 1	1, 5-10 / 1	1-5 / 1

Table 10: Retrieval range for the 4D-Lung and CVC-ClinicDB dataset. **Blue** indicates the retrieval range for positive-point prompts. **Red** indicates the retrieval range for negative-point prompts.

Dataset	<i>Meta</i>		<i>LoRA</i>		<i>Full-Fine-Tune</i>	
	<i>huge</i>		<i>base</i>	<i>large</i>	<i>base</i>	<i>large</i>
4D-Lung	1-2 / 45		1-3 / 1	1-3 / 1	1-3 / 1	1-3 / 1
CVC-ClinicDB	1-5 / 1-3		1-3 / 1-3	1-2 / 1-3	1-2 / 1	1-5 / 1-3

In this section, for reproducibility, we provide the details of the retrieval range during the test time for the COCO-20ⁱ (Nguyen & Todorovic, 2019), FSS-1000 (Li et al., 2020), LVIS-92ⁱ (Liu et al., 2023), and Perseg (Zhang et al., 2023) dataset in Table 9, the 4D-Lung (Hugo et al., 2016) and CVC-ClinicDB (Bernal et al., 2015) dataset in Table 10.

The final number of positive-point and negative-point prompts is determined by our distribution-guided retrieval approach. Below, we explain how the retrieval range is determined in Table 10. For *LoRA* and *Full-Fine-Tune*, the retrieval range is determined based on the validation set of the *internal* datasets. We uniformly sample positive-point and negative-point prompts on the ground-truth mask and perform interactive segmentation. The number of prompts is increased until the improvement becomes marginal, at which point this maximum number is set as the retrieval range for *external* test datasets. On the 4D-Lung dataset, we consistently set the number of negative-point prompts to 1 for these two types of models. This decision is informed by conclusions from previous works (Ma et al., 2024a; Huang et al., 2024), which suggest that the background and semantic target can appear very similar in CT images, and using too many negative-point prompts may confuse the model. On the CVC-ClinicDB dataset, the endoscopy video is in RGB space, resulting in a relatively small domain gap (Matsoukas et al., 2022) compared to SAM’s pre-trained dataset. Therefore, for *Meta*, we use the same retrieval range as the *Full-Fine-Tune* large model. In contrast, on the 4D-Lung dataset, CT images are in grayscale, leading to a significant domain gap (Matsoukas et al., 2022) compared to SAM’s pre-trained dataset. Consequently, we set the retrieval range for positive-point prompts to 2 to avoid outliers and fixed the number of negative-point prompts to a large constant (*i.e.*, 45) rather than a range, to ensure the model focuses on the semantic target. These values were not further tuned.

D Additional Visualization

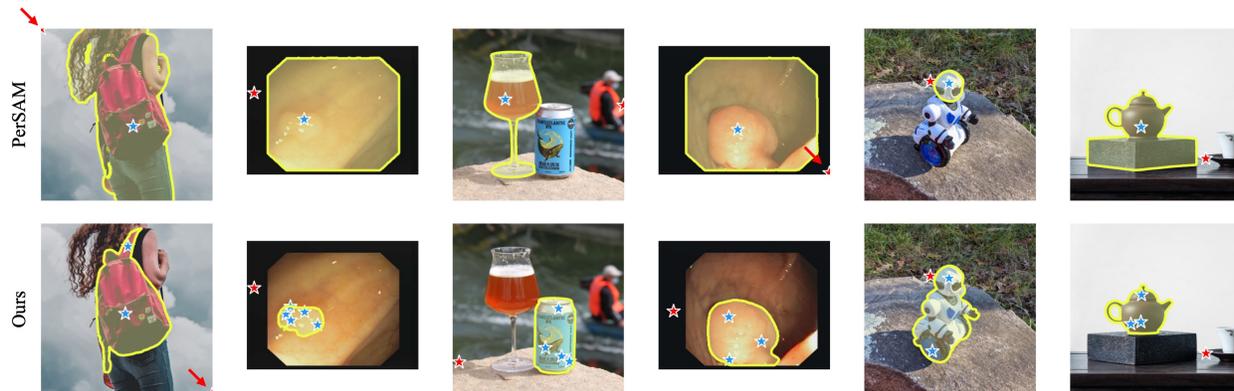


Figure 12: Additional qualitative results: (Columns 1–4) Full images from earlier illustrations; (Columns 5–6) Additional comparisons with PerSAM. Note that the negative-point prompt can sometimes differ between P²SAM and PerSAM, as the similarity matrix changes when using part-level features.

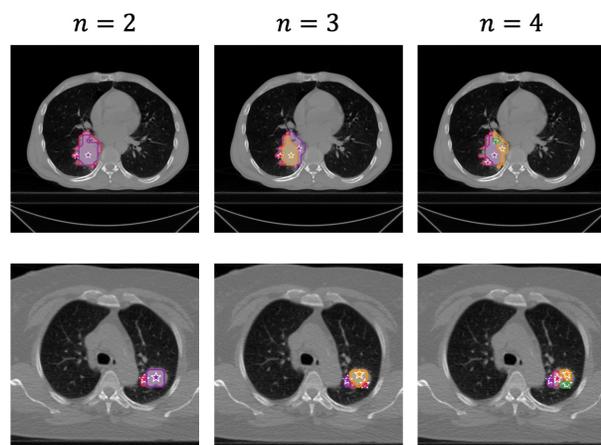


Figure 13: Visualization results on the 4D-Lung dataset, based on a varying number of part-level features.

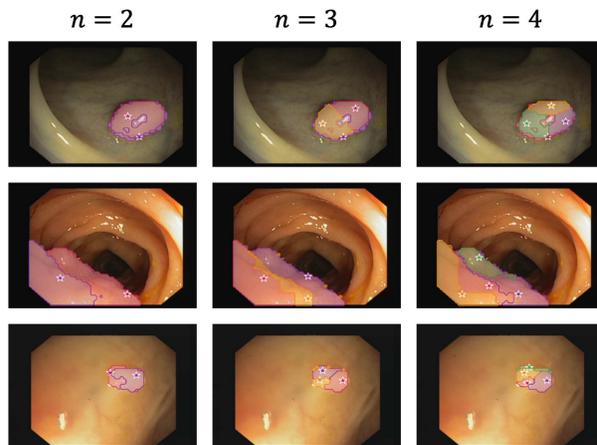


Figure 14: Visualization results on the CVC-ClinicDB dataset, based on a varying number of part-level features.

In this section, we first provide the full images in Figure 12 that were presented in Section 1 to eliminate any possible confusion. Then, to provide deeper insight into our part-aware prompt mechanism and distribution-guided retrieval approach, we present additional visualization results on the 4D-Lung (Hugo et al., 2016) dataset, the CVC-ClinicDB (Bernal et al., 2015) dataset, and the PerSeg (Zhang et al., 2023) dataset. These visualizations are based on a varying number of part-level features, offering a clearer understanding of how the part-aware prompt mechanism adapts to different segmentation tasks and domains. In Figure 13 and 14, we observe that an appropriate number of part-level features can effectively divide the tumor into distinct parts, such as the body and edges for non-small cell lung cancer, and the body and light point (caused by the camera) for the polyp. This illustrates how P²SAM can assist in cases of incomplete segmentation. In Figure 15, we observe that an appropriate number of part-level features can effectively divide the object into meaningful components, such as the pictures, characters, and aluminum material of a can; the legs and platforms of a table; or the face, ears, and body of a dog. These parts can merge naturally based on texture features when using the appropriate number of part-level features, whereas using too many features may result in over-segmentation. Our retrieval approach, on the other hand, helps determine the optimal number of part-level features for each specific case.

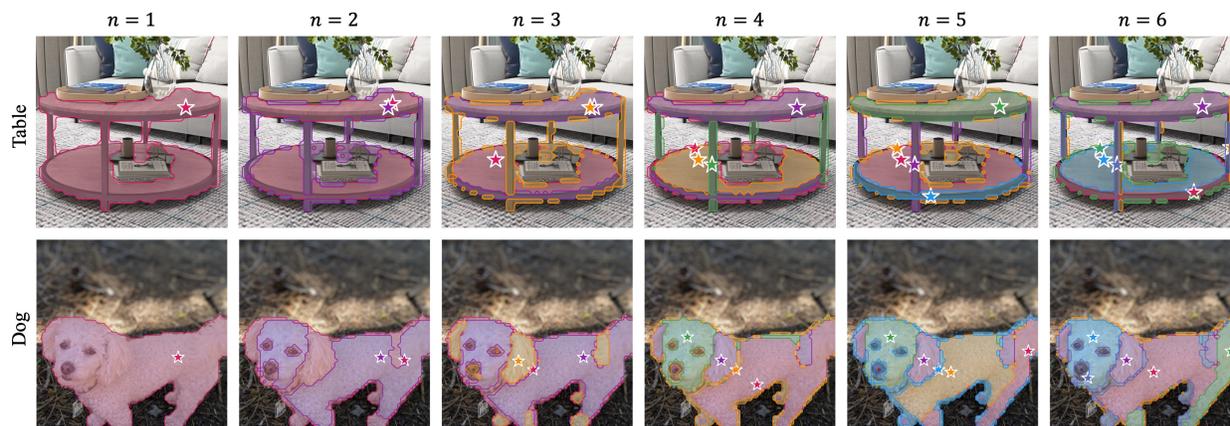


Figure 15: Visualization results on the PerSeg dataset, based on a varying number of part-level features.

E Discussion

Table 11: Results of *direct-transfer* on CVC-ClinicDB. The model is trained on Kvasir-SEG with different pre-training weights.

Med-SAM	83.85
SAM	84.62

Table 13: Comparison with GF-SAM on the CVC-ClinicDB dataset. \star indicates using DI-NOv2 for a better performance.

Method	<i>Meta</i>	<i>Full-Fine-Tune</i>
GF-SAM	60.55 \star	87.57
Matcher	63.54	87.15
P ² SAM	66.45	88.76

Table 12: Results of interactive segmentation on the internal Kvasir-SEG validation dataset and the external CVC-ClinicDB dataset. We use *Full-Fine-Tune large*^{312.5M} here.

Dataset	No Prompt	P ² SAM Point Prompt	1 Positive Prompt	Box
Kvasir-SEG	93.27	-	95.15	95.57
CVC-ClinicDB	86.68	88.76	88.99	92.35

Table 14: Results of one-shot part segmentation on the PASCAL-Part dataset. Note that all methods utilize SAM’s encoder for fairness.

Method	animals	indoor	person	vehicles	mean
Matcher	29.29	56.30	21.04	37.02	33.66
PerSAM	19.9	51.8	18.6	32.0	30.1
P ² SAM	20.29	54.82	19.62	34.21	32.24
P ² SAM <i>w. neg</i>	20.54	54.65	20.59	36.91	33.17

Baseline Results. In this paper, we treat MedSAM (Ma et al., 2024a) with a human-given box prompt as the baseline for the 4D-Lung dataset (Hugo et al., 2016), DuckNet (Dumitru et al., 2023) as the baseline for the CVC-ClinicDB dataset (Bernal et al., 2015). We acknowledge that MedSAM is widely used as a baseline across many benchmarks (Antonelli et al., 2022; Ji et al., 2022). However, these comparisons primarily focus on internal validation. MedSAM has the potential to outperform many models on external validation sets due to its pre-training on a large-scale medical image dataset. While there is no direct evidence to confirm this, DuckNet (Dumitru et al., 2023) suggests that large-scale pre-trained models generally outperform others on external validation sets, even if they lag behind on internal validation. Among studies (Butoi et al., 2023; Wong et al., 2023; Ma et al., 2024a;b; Wu & Xu, 2024) that aim to develop promptable segmentation models specifically for medical image segmentation, UniverSeg (Butoi et al., 2023)’s performance may decline significantly with only one-shot support set, and both ScribblePrompt (Wong et al., 2023) and One-Prompt (Wu & Xu, 2024) are trained on much smaller datasets. As we focus on segmenting external patient samples that lie outside the training distribution in a one-shot manner. Therefore, we argue that the model’s generalization ability is critical for achieving superior performance. The 4D-Lung dataset (Hugo et al., 2016) is a relatively

new benchmark for longitudinal data analysis, and no standard benchmark for comparison was available at the time this work was conducted. In addition, during evaluation, we supplemented MedSAM with a human-given box prompt, making it a very fair baseline for this work.

Baseline Methods. In this paper, we treat SAM-based methods such as PerSAM (Zhang et al., 2023) and Matcher (Liu et al., 2023) as our primary baselines and also compare with PANet (Wang et al., 2019b). We do not include other backbone methods like ScribblePrompt (Wong et al., 2023) and One-Prompt (Wu & Xu, 2024) because they primarily focus on interactive segmentation, just similar to MedSAM (Ma et al., 2024a), which is the baseline we compare in Table 3. On the other hand, utilizing other prompt modalities, such as scribble, mask, and box, presents challenges for solving the patient-adaptive segmentation problem, as it is difficult to represent prior data in these formats. In this work, we adopt a more flexible prompt modality: point prompts. Although it may be possible to convert our multiple-point prompts into a scribble prompt by connecting them together, we leave the exploration of this direction for future work. Consequently, the most relevant baseline methods remain SAM-base methods like PerSAM and Matcher.

Here, we evaluate a more recent SAM-base method, GF-SAM (Zhang et al., 2025). Similar to Matcher, GF-SAM utilizes DINOv2 to extract patch-level features; however, GF-SAM is a hyper-parameter-free method based on graph analysis. In Table 13, we evaluate GF-SAM on the CVC-ClinicDB dataset (Bernal et al., 2015) using both a natural image pre-trained encoder (*Meta*) and a medically adapted encoder (*Full-Fine-Tune*). With the natural image pre-trained encoder, P²SAM outperforms both GF-SAM and Matcher, since patch-level features are less robust than part-level features when there is domain gap between pre-training data and test data. However, GF-SAM fails to surpass Matcher in this task, which contrasts with its superior performance on natural image segmentation tasks. We hypothesize that this is because GF-SAM is a hyper-parameter-free method, and factors such as the number of point prompts, the number of clusters, and the threshold value may be more sensitive when there is a domain gap between the pre-training data and the test data. GF-SAM outperforms Matcher with the medically adapted encoder, but still lags behind P²SAM, as the encoder is adapted for medical segmentation tasks and still lacks patch-level objectives. This result, along with the findings in Table 8—where P²SAM with *base* SAM outperforms PerSAM with *huge* SAM by 0.7% mIoU and *base* SAM by 26.0% mIoU on the PerSeg dataset—further underscores that P²SAM is a more robust method when the model exhibits weaker representations, a scenario more prevalent in medical image analysis.

Pre-trained Model. In this work, we choose to adapt SAM to the medical image domain using the SA-1B pre-trained model weights rather than weights from MedSAM for two reasons. First, although MedSAM fine-tunes SAM (SA-1B pre-trained) on a large-scale medical segmentation dataset, its fine-tuning dataset is still 1,000 times smaller than SAM’s pre-training dataset (1M vs. 1B). Since model generality after adaptation is crucial for our work, we assume that SAM remains a better starting point, despite MedSAM being a strong option for zero-shot medical segmentation. Second, MedSAM only provides the SAM-Base pre-trained model, whereas our results in Table 1 and Table 2 demonstrate that larger models (i.e., *large*) can further enhance performance across various tasks. In Table 11, we provide the *direct-transfer* result on the CVC-ClinicDB dataset, the model is trained on the Kvasir-SEG dataset with Med-SAM pre-trained weights and SA-1B pre-trained weights. The result follows our assumption and the discussion in MedSAM (Ma et al., 2024a) and its successor (Ma et al., 2024b), that with a specific task, maybe fine-tune from SAM is still a better choice.

Interactive Segmentation. As mentioned in Section 3.3 and detailed in Appendix B, we closely adhere to SAM’s interactive training strategy when adapting it with medical datasets. Therefore, our medically adapted model retains its interactive segmentation capability. In Table 12, we present both internal evaluation results on the Kvasir-SEG dataset’s validation set and external evaluation results on the CVC-ClinicDB dataset. First, as discussed in Section 4.2 and Appendix B, since we have a specific segmentation target, our adapted model does not need to be ambiguity-aware, allowing a human-given single positive-point prompt to achieve good performance. P²SAM lags only slightly behind this result while operating fully automatically. For the human-given box prompt, it is not surprising that it outperforms P²SAM, as a box prompt is a strong prompt that essentially requires the provider to know the lesion’s location.

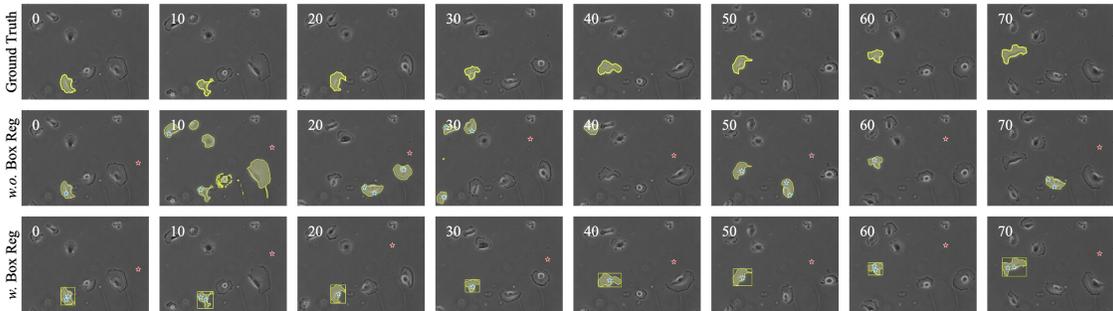


Figure 16: Qualitative results of single-cell segmentation on the PhC-C2DH-U373 dataset. The second row highlights the challenge P²SAM faces in handling multiple similar objects. The third row demonstrates that P²SAM can overcome this challenge with a cost-free regularization.

Part Segmentation. We acknowledge that P²SAM’s design was not initially focused on part segmentation but on enhancing the medical image segmentation model’s generality by providing more precise and informative prompts. We conduct the part segmentation task on the PASCAL-Part dataset (Everingham et al., 2010). Note that all methods use SAM (*Meta*) as the backbone model. Part segmentation with SAM typically relies more on additional prompt modalities, such as box prompts, or diverse mask candidates. For example, Matcher employs a random point-prompt sampling strategy to make their proposed mask candidates more diverse, potentially slowing down the algorithm. In Table 14, when compared with PerSAM, P²SAM consistently shows benefits (i.e., +2.23% mIoU). However, P²SAM is surpassed by Matcher (i.e., -1.42% mIoU). For P²SAM it is reasonable to provide additional negative-point prompts in part segmentation task because a portion of the background is correlated between the reference and target images (i.e., both refer to the rest of the object). Therefore, we additionally provide negative-point prompts to P²SAM (P²SAM *w. neg*), which further improves segmentation performance (i.e., +0.93% mIoU) and brings P²SAM on par with Matcher. While achieve slightly better performance, Matcher utilizes 128 sampling iterations for the part segmentation task, making it much slower (x3) than both PerSAM and P²SAM.

Similar Objects. P²SAM demonstrates improvements in the backbone’s generalization across domain, task, and model levels. At the task level, we have already shown how P²SAM enhances performance for NSCLC segmentation in patient-adaptive radiation therapy and polyp segmentation in endoscopy videos. However, when addressing specific tasks that involve multiple similar targets, P²SAM may fail due to the lack of instance-level objective. Although this scenario is uncommon in patient-adaptive segmentation, we acknowledge that P²SAM faces the same challenge of handling multiple similar objects as other methods (Zhang et al., 2023; Liu et al., 2023). In Figure 16, we present an example of single-cell segmentation on the PhC-C2DH-U373 dataset (Maška et al., 2014), which goes beyond the patient-specific setting. In Figure 16, the second row illustrates that P²SAM fails to segment the target cell due to the presence of many similar cells in the field of view. However, given the slow movement of the cell, we can leverage its previous information to regularize the current part-aware prompt mechanism. The third row in Figure 16 demonstrates that when using the bounding box from the last frame, originally propagated from the reference frame, to regularize the part-aware prompt mechanism in the current frame, P²SAM achieves strong performance on the same task. Since the bounding box for the first frame can be generated from the ground truth mask, which is already available, this regularization incurs no additional cost. Utilizing such tailored regularization incorporating various prompt modalities, we showcase our approach’s flexible applicability to other applications.

F Equations

In this section, we provide details on the equation mentioned in Section 3.2.

Wasserstein Distance. In Equation 3, we use $\mathcal{D}_w(\cdot, \cdot)$ to represent the Wasserstein Distance. Here we provide the details of this function. Suppose that features in the reference image $F_R \in \mathbb{R}^{n_r \times d}$ and features in the target image $F_T \in \mathbb{R}^{n_t \times d}$ come from two discrete distributions, $F_R \in \mathbf{P}(\mathbb{F}_{\mathbb{R}})$ and $F_T \in \mathbf{P}(\mathbb{F}_{\mathbb{R}})$, where

$F_R = \sum_{i=1}^{n_r} u_i \delta_{f_r}^i$ and $F_T = \sum_{j=1}^{n_t} v_j \delta_{f_t}^j$; δ_{f_r} being the Delta-Dirac function centered on f_r and δ_{f_t} being the Delta-Dirac function centered on f_t . Since F_R and F_T are both probability distributions, sum of weight vectors is 1, $\sum_i u_i = 1 = \sum_j v_j$. The Wasserstein distance between F_R and F_T is defined as:

$$\mathcal{D}_w(F_R, F_T) = \min_{\mathbf{T} \in \Pi(u, v)} \sum_i \sum_j \mathbf{T}_{ij} \cdot \frac{F_R^i \cdot F_T^j}{\|F_R^i\|_2 \cdot \|F_T^j\|_2} \quad (5)$$

where $\Pi(u, v) = \{\mathbf{T} \in \mathbb{R}_+^{n_r \times n_t} | \mathbf{T} \mathbf{1}_m = u, \mathbf{T}^\top \mathbf{1}_n = v\}$, and \mathbf{T} is the transport plan, interpreting the amount of mass shifted from F_R^i to F_T^j .