# Exploring Transfer Learning in Medical Image Segmentation using Vision-Language Models

**Kanchan Poudel**[*]  **Manish Dhakal**[*]  **Prasiddha Bhandari**[*]  **Rabin Adhikari**[*]
**Safal Thapaliya**[*]      **Bishesh Khanal**

Nepal Applied Mathematics and Informatics Institute for research (NAAMII)
{kanchan.poudel,manish.dhakal,prasiddha.bhandari,rabin.adhikari,
safal.thapaliya,bishesh.khanal}@naamii.org.np

## Abstract

Medical image segmentation with deep learning is an important and widely studied topic because segmentation enables quantifying target structure size and shape that can help in disease diagnosis, prognosis, surgery planning, and understanding. Recent advances in the foundation Vision-Language Models (VLMs) and their adaptation to segmentation tasks in natural images with Vision-Language Segmentation Models (VLSMs) have opened up a unique opportunity to build potentially powerful segmentation models for medical images that enable providing helpful information via language prompt as input, leverage the extensive range of other medical imaging datasets by pooled dataset training, adapt to new classes, and be robust against out-of-distribution data with human-in-the-loop prompting during inference. Although transfer learning from natural to medical images for image-only segmentation models has been studied, no studies have analyzed how the joint representation of vision-language transfers to medical images in segmentation problems and understand gaps in leveraging their full potential.

We present the first benchmark study on transfer learning of VLSMs to 2D medical images with thoughtfully collected 11 existing 2D medical image datasets of diverse modalities with carefully presented 9 types of language prompts from 14 attributes. Our results indicate that VLSMs trained in natural image-text pairs transfer reasonably to the medical domain in zero-shot settings when prompted appropriately for non-radiology photographic modalities; when finetuned, they obtain comparable performance to conventional architectures, even in X-rays and ultrasound modalities. However, the additional benefit of language prompts during finetuning may be limited, with image features playing a more dominant role; they can better handle training on pooled datasets combining diverse modalities and are potentially more robust to domain shift than the conventional segmentation models. The code and datasets are released at `https://github.com/naamiinepal/med vlsm`.

## 1  Introduction

Medical image segmentation is essential in clinical workflow for diverse applications such as computer-aided diagnosis, prognosis, surgery planning, or population-based studies. The state-of-the-art supervised segmentation models, following similar encoder-decoder architecture [43] with Convolutional Neural Networks (CNNs) or Vision Transformers (ViTs) [10], have shown promising results in well-curated datasets across almost all imaging modalities, anatomies, and diseases [5, 19, 20, 23, 37, 39, 54, 61]. However, each model, trained on a particular curated dataset (or

---

[*]Equal Contribution. The order is in the ascending order of the authors' first names.

finetuned when the encoder backbone is pretrained on ImageNet [9]), works only for the predefined set of foreground classes on the specific image modality and anatomy, and is primarily unable to leverage potentially useful auxiliary information (other than the input image) that could build robustness against domain shift and out-of-distribution data. Thus, in applications except for extensive population-based studies, clinicians mostly resort to manual or semi-automated interactive segmentation approaches. Improving the speed and quality of the interactive segmentation can substantially impact the clinical workflow and healthcare service delivery.

Recent introduction of Vision-Language Models (VLMs) [14, 22, 25, 33, 41, 48, 57] and their adaptation to segmentation tasks for natural images with Vision-Language Segmentation Models (VLSMs) [36, 42, 52] have opened up an opportunity to learn powerful text-image joint representation and provide auxiliary information as input via language prompts during image segmentation. Compared to visual prompts such as points or boxes, language prompts during inference are, by design, more interpretable or explainable and can provide complex auxiliary information regarding texture, shape, and spatial relationships of normal and abnormal structures. Foundation VLSMs trained on large-scale image-text pairs and the ability to inject auxiliary prompts via language prompt could enable powerful medical image segmentation models that can leverage the separately curated medical datasets (often of small or medium size) by pooled dataset training, adapt to new classes without changing the architecture, provide more explainable results, and potentially have robust human-in-the-loop segmentation systems for out-of-distribution data and domain adaptation.

Although transfer learning from natural images to medical images for image-only representation learning has been studied and widely used [2, 8, 15], the use or study of transferring joint vision-language representation is at a nascent stage. Qin et al. [40] studied transfer learning for one VLM designed for object detection task, showing promising zero-shot results and surpassing baseline supervised model when finetuning the VLMs with suitable language prompts on small-sized medical images. However, the two critical questions remain unanswered: (**i**) How well this holds for the segmentation task across multiple VLSMs is unclear. Segmentation is a more difficult task than object detection and arguably more important and ubiquitous in medical imaging because accurate segmentation can provide quantitative measures of target structures. (**ii**) There is a lack of understanding of the nuanced role of the language prompt vs. image during finetuning and the ability of the VLSMs to handle pooled dataset training and out-of-distribution data.

We propose the first benchmark study on transfer learning of VLSMs based on two VLMs – one pretrained on natural image-text pairs (CLIP [41]) and another pretrained on medical image-text pairs (BiomedCLIP [59]). In addition to CLIPSeg [36] and CRIS [52] (the two VLSMs proposed in the literature with different architecture designs and trained on a large dataset of natural image segmentation data to learn pixel-token level representation), we build two more VLSMs from BiomedCLIP: (**i**) BiomedCLIPSeg-D – obtained by adding pretrained CLIPSeg decoder to BiomedCLIP, and (**ii**) BiomedCLIPSeg – obtained by initializing the decoder from scratch.

Our significant contributions are listed below.

1. We carefully select a wide range of datasets (11) covering four common 2D medical image modalities and six target structures comprising anatomical structures and pathology, such as tumors, with binary and multi-class segmentation tasks.

2. In addition to collecting existing datasets, we enrich them for studying and benchmarking VLSMs with a rich set of language prompts. The prompts are generated from our proposed automated and practical method, which uses diverse sources such as image metadata, VQA model, and segmentation masks and provides valuable information at individual image samples and target class levels.

3. We conduct extensive experiments using four VLSMs with diverse datasets and carefully designed prompts to explore the nuanced relationship between language and image when adapting the joint representation for medical images with limited data, study robustness against domain shift, ability to handle pooled dataset comprising of diverse modalities, attributes, and target masks.

4. We make our framework, source code, and the generated language prompts open-source, fostering transparency and reproducibility.

## 2 Image Segmentation using Foundation Vision Language Models

### 2.1 Vision Language Pretraining and Foundation Models

Recently proposed foundation VLMs like Contrastive Language-Image Pretraining (CLIP) [41] jointly train a transformer-based text encoder and an image encoder on large-scale text-image pairs. They use a contrastive loss to maximize the similarity of the image and text embeddings of correct pairs while minimizing the embedding similarity of incorrect pairings. Instead of contrastive learning, FLAVA [48] adds another multimodal encoder that inputs the embeddings from individual encoders and trains the model using different loss functions corresponding to vision, language, and multi-model tasks. Other VLMs, which are similar to CLIP, attempt to address some of the issues, such as the requirement of a large number of image-language pairs [33], noisy image-text pairs [25, 32], and high computational complexity [57].

Foundation VLMs for medical imaging applications in the literature primarily build upon the VLMs developed with natural images using one of the two approaches: (**1**) finetune a pretrained VLM with medical image-text pairs (Seibold et al. [46], PubMedCLIP [12]), or (**2**) pretrain the VLM of the same architecture from scratch with image-text pairs (ConVIRT [60], MedCLIP [53], MedKLIP [56], BiomedCLIP [59]).

Although these foundation models are intended to be helpful for a wide range of downstream tasks, their general approach to learning an embedding that aligns the whole image to the entire input text with global embedding is sub-optimal for dense prediction tasks like segmentation.

### 2.2 Vision Language Segmentation Models (VLSMs) with Pixel-Token alignment

Segmentation tasks may benefit more from explicitly aligning images and text descriptions. Recent state-of-the-art VLSMs extend CLIP to segmentation tasks by adding a decoder trained to produce a segmentation map from CLIP's vision and language embeddings. DenseCLIP [42] has proposed vision language decoders on top of the CLIP encoders that use the pixel-text score maps of the limited class prompts. In contrast to DenseCLIP, CLIPSeg [36] and CRIS [52] enforce zero-shot segmentations by giving the output as pixel-level activations for the given text or image prompt. ZegCLIP [62] has introduced tune prompts and associated the image information to the text encodings before patch-text contrasting to reduce the seen classes' overfit.

Although there are some specific architectures to learn joint embeddings from image and text prompts for particular datasets, such as TGANet [50] for polyps in endoscopy images, to our knowledge, there are no VLSMs well studied for medical images.

## 3 Method

### 3.1 CLIP- and BiomedCLIP-based Medical VLSMs

We develop four medical VLSMs using two different strategies: (**i**) Finetuning two state-of-the-art CLIP-based VLSMs, **CLIPSeg** [36] and **CRIS** [52], that are pretrained on natural image-text pairs, and (**ii**) Building two new VLSMs specific to medical domain by adding a decoder to the foundation VLM BiomedCLIP [59], pretrained on medical image-text pairs. The two new VLSMs are **BiomedCLIPSeg-D** – obtained by adding pretrained CLIPSeg decoder to BiomedCLIP, and **BiomedCLIPSeg** – obtained by randomly initializing the decoder. These four models are trained on the triplet of medical images, segmentation masks, and text prompts. To compare how VLMs pretrained on a large corpus of medical image performs compared to the one pretrained on open-domain VLMs, we chose BiomedCLIP as the base VLM because it was trained on diverse medical imaging modality and the same CLIP architecture.

Figure 1 shows the general architecture of our VLSMs where the embeddings learned during pretraining from the text and image encoders (either of CLIP [41] or BiomedCLIP [59]) are aggregated and fed to vision language decoder to generate the binary segmentation mask. We study CLIPSeg and CRIS models for both the zero-shot and finetuning settings, but only the later setting for BiomedCLIPSeg and BiomedCLIPSeg-D as they do not have an end-to-end pretrained encoder-decoder.
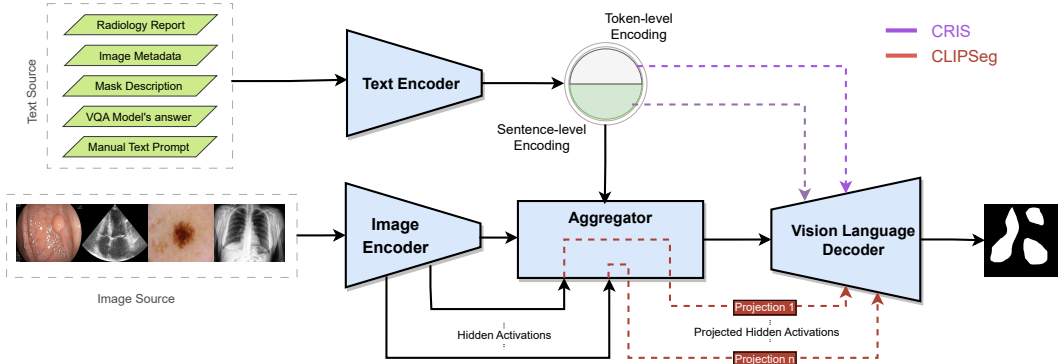
Figure 1: **The basic architecture of CRIS and CLIPSeg/BiomedCLIPSeg* VLSMs.** The VLSM consists of a Text Encoder, an Image Encoder, a Vision-Language Decoder (VLD), and an Aggregator. The input to the model is a pair of an image and a language prompt fed to the Image Encoder and Text Encoder, respectively. The Aggregator generates intermediate representations utilizing image-level, sentence-level, or word-level representations to feed to the VLD.

CLIPSeg and CRIS[2] were trained on PhraseCut [55] with $340,000$ text-image pairs and RefCOCO [28] with $142,210$ text-image pairs, respectively. Architecture-wise, CLIPSeg has the text embedding at the neck of the vision encoder-decoder network. At the same time, CRIS gives more focus to the text branch by inserting the word-level features at the neck and sentence-level features at the penultimate layer of the decoder. CRIS's adaptation of the Image Encoder to insert word-level features supports only the pretrained CNN-based CLIP backbone but not the ViT-based CLIP backbone. On the contrary, CLIPSeg works with both CNN and ViT backbones; BiomedCLIP pretrained CLIP architecture with ViT-based backbone on $15$ million image-text pairs extracted from biomedical research articles in PubMed Central (PMC-15M) [59].

## 3.2 Datasets

We collected 11 2D medical imaging datasets of diverse modalities, organs, and pathologies covering both radiology and non-radiology images for binary and multi-class segmentation tasks (see Table 1). All the datasets are used for finetuning separately or combined (as a single pooled dataset) except the last three endoscopy datasets (ETIS, ColonDB, and CVC300), which are used only as the test split to study domain shift robustness.

Table 1: An overview of our datasets compared across the dimensions of category, modality, organ, foreground classes, and their splits. The datasets with multiple foreground classes signify the multi-class segmentation tasks.

| Category | Modality | Organ | Name | Foreground Class(es)/Ground Truth mask name | # train/val/test |
|---|---|---|---|---|---|
| Non-Radiology | Endoscopy | Colon | Kvasir-SEG [24] | Polyp | 800/100/100 |
| | | | ClinicDB [4] | | 490/61/61 |
| | | | BKAI [3, 38] | | 800/100/100 |
| | | | ETIS [47] | | 0/0/196 |
| | | | ColonDB [49] | | 0/0/380 |
| | | | CVC300 [51] | | 0/0/60 |
| | Photography | Skin | ISIC 2016 [17] | Skin Lesion | 810/90/379 |
| | | Foot | DFU 2022 [29] | Foot Ulcer | 1600/200/200 |
| Radiology | Ultrasound | Heart | CAMUS [31] | Myocardium, Left ventricular, and Left atrium cavity | 4800/600/600 |
| | | Breast | BUSI [1] | Benign and Malignant Tumors | 624/78/78 |
| | X-Ray | Chest | CheXlocalize [45] | Atelectasis, Cardiomegaly, Consolidation, Edema, Enlarged Cardiomediastinum, Lung Lesion, Lung Opacity, Pleural Effusion, Pneumothorax, and Support Devices | 1279/446/452 |

---

[2]We used unofficial weights from the GitHub issue (`https://github.com/DerrickWang005/CRIS.pytorch/issues/3`) since the authors haven't released the model weights yet.

### 3.3 Generating Language Prompts

Language prompt provides powerful flexibility to describe images and provide context, potentially creating a more robust and rich interactive segmentation. However, manually entering individual image-specific prompts is impractical for a large-scale evaluation. Thus, we create automatic prompts for large-scale evaluations of medical VLSMs, using common VLSM concepts like size, position, and color, specific medical concepts like gender, age, and pathology, and add manual prompts only for general class-level information common for all samples in a dataset. Diverse sources were considered for extracting useful information, including VQA, medical image metadata, and automated image processing. The attributes and information used to generate language prompts are summarized below:

- *Number*, *size* and *relative location of the target on image* are generated by applying image processing on segmentation masks, as suggested by Tomar et al. [50], for non-radiology images and whichever of the three seemed relevant to radiology image datasets.
- Motivated from Qin et al. [40], we query VQA for *shape* and *color* in non-radiology datasets, and only *shape* in radiology datasets since *color* is irrelevant for grayscale radiology images.
- *General class information* for photographic images from online medical journals [3]. Qin et al. [40] used PubMedBERT [16], an MLM trained in medical texts, to extract general class information. However, our experiments with this model gave unreliable outputs [4], driving us to get this information manually from online medical journals.
- *Medical reports* and attributes extracted from other image-specific metadata, like *age* and *gender* of patients, *image quality*, *cardiac cycle*, *tumor type*, whenever available.

Table 2 shows how 14 attributes relevant to various datasets (**a1**-**a14**) are combined to provide nine different prompt types (**P1**-**P9**), with additional prompt **P0** for empty prompt. We ordered the attributes in the prompt based on their perceived importance, starting with the class name as the essential attribute to inform models of what target structure to segment. While many more language prompts exist when permuting the order of attributes order, it requires enormous resources. Instead, we focused on designing experiments to analyze how well the concepts were learned by the VLSMs, such as evaluating the performance when intentionally giving wrong attribute values or randomly mapping to uncommon English words without semantic meaning.

Table 2: Different prompts are formed for each dataset using combinations of 14 potential attributes. Although some attributes, like *Pathology*, are specific to some particular datasets, others, like *Class Keywords*, are common to all the datasets.

**Attributes → a1:** Class Keyword; **a2:** Shape; **a3:** Color; **a4:** Size; **a5:** Number; **a6:** Location; **a7:** General Class Info; **a8:** View; **a9:** Pathology; **10:** Cardiac Cycle; **a11:** Gender; **a12:** Age; **a13:** Image Quality; **a14:** Tumor Type

| Prompts → / Datasets ↓ | P1 | P2 | P3 | P4 | P5 | P6 | P7 | P8 | P9 |
|---|---|---|---|---|---|---|---|---|---|
| **Non-Radiology** | a1 | a1a2 | a1a2a3 | a1a2a3a4 | a1a2a3a4a5 | a1a2a3a4a6 | a1a7 | a1a2a3a4a5a7 | a1a2a3a4a5a6a7 |
| Example Prompt | **P9 →** one small pink round polyp which is **often a bumpy flesh in rectum** located in **center** of the image | | | | | | | | |
| **CheXlocalize** | a1 | a1a8 | a1a2a8 | a1a2a6a8 | a1a2a6a8a9 | a1a9 | N/A | N/A | N/A |
| Example Prompt | **P5 → Airspace Opacity** of shape **rectangle**, and located in **right** of the **frontal** view of a Chest Xray. **Enlarged Cardiomediastinum, Cardiomegaly, Lung Opacity, Consolidation, Atelectasis, Pleural Effusion** are present. | | | | | | | | |
| **CAMUS** | a1 | a1a8 | a1a8a10 | a1a8a10a11 | a1a8a10a11a12 | a1a8a10a11a12a13 | a1a8a10a11a12a13a2 | N/A | N/A |
| Example Prompt | **P7 → Left ventricular cavity** of **triangular shape** in **two-chamber view** in the cardiac ultrasound at the end of the **diastole cycle** of a **40-year-old female** with **poor image quality**. | | | | | | | | |
| **BUSI** | a1 | a1a14 | a1a14a5 | a1a14a5a4 | a1a14a5a4a6 | a1a14a5a4a6a2 | N/A | N/A | N/A |
| Example Prompt | **P6 → Two medium square-shaped benign tumors** at the **center, left** in the breast ultrasound image. | | | | | | | | |

### 3.4 Implementation Details

For the VLSMs, we have kept the experimental settings close to the original while maintaining consistency with minimal hyperparameter changes. The models are finetuned and inferred in NVIDIA GeForce RTX 3090 and NVIDIA Titan Xp GPUs with a batch size of 16.

We use Adam [30] optimizer with a constant learning rate of $10^{-3}$ with no weight decay with the weighted sum of Dice and Binary Cross Entropy (BCE) losses with weights of 1 and 0.2, respectively

---

[3]More in Table 9 in the supplementary section.

[4]See Table 5 in the supplementary section

for CLIPSeg. Similarly, following the original implementation, for CRIS, we use the same optimizer with the learning rate of $10^{-6}$ for the encoders and $10^{-5}$ for the decoder, a dropout of $0.2$ for the decoder and BCE as the loss function.

CLIPSeg and CRIS internally resize the three-channeled input images to $416 \times 416$ and $352 \times 352$, respectively. The dice scores mentioned in the paper are calculated after resizing the output of the models back to the original size (before respective resizing). We normalize the resized images with means and standard deviations provided by the respective models and haven't performed other preprocessing and post-processing to access the models' raw performance.

## 4    Results

This section first presents experimental results in zero-shot settings (for CRIS and CLIPSeg) and finetuned settings (for all four VLSMs) using a maximum of nine prompts on all the datasets. This is followed by a more subtle look into how well the VLSMs capture concepts represented by different attributes and the influence on segmentation output when wrong information is provided. Finally, the robustness of the VLSMs to handle diverse datasets and comparison against common segmentation models are reported.

**VLSMs adapt better to non-radiology images in zero-shot settings.**    Figure 2 shows that both CRIS and CLIPSeg barely work in zero-shot settings for radiology datasets except for CRIS in breast ultrasound (BUSI) but get a Dice score (DSC) in the range of $20 - 60$ for endoscopy datasets. The highest DSC is for the skin dataset (ISIC), reaching up to $67.98$ for CLIPSeg. While adding more attributes to the prompt seems to have generally increased the performance, the gain is inconsistent across all prompts and datasets.
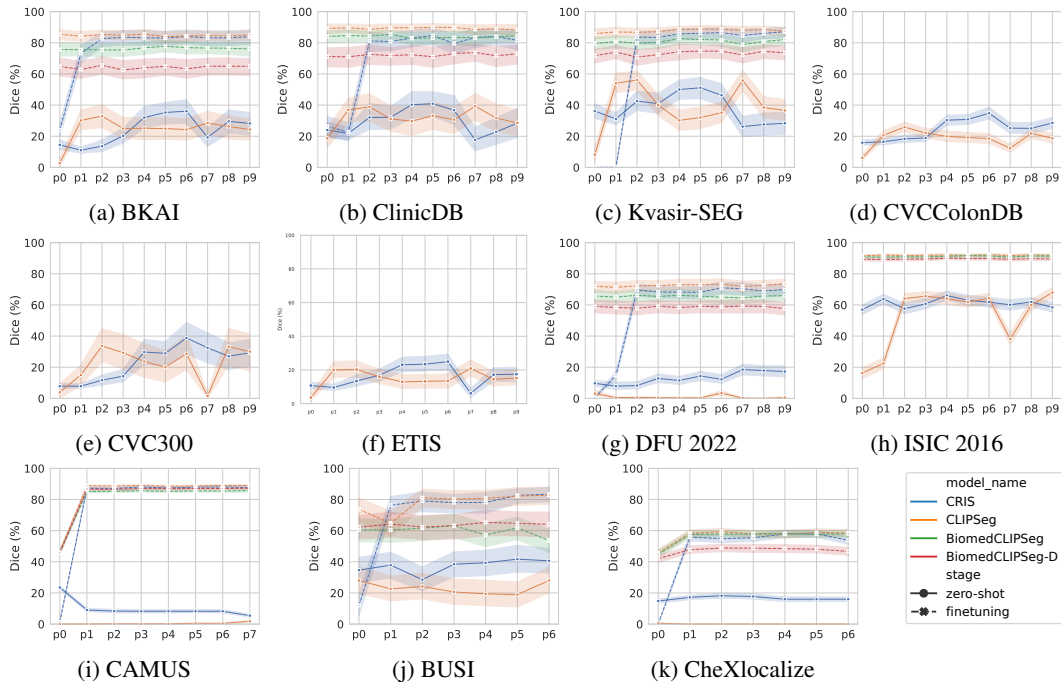


Figure 2: Zero-shot and finetuning performance of CRIS, CLIPSeg, BiomedCLIPSeg, and BiomedCLIPSeg-D model on non-radiology ( first two rows) and radiology datasets (last row). Finetuning using the prompts shows performance improvement compared to the empty prompt, which is more prominent in multi-class settings. However, using the label name and adding additional prompts does not significantly affect the model performance.

**Image-specific-attributes or general descriptions?**    In the zero-shot setting, CRIS has better segmentation performance in almost all endoscopy datasets when the prompt contains multiple image-specific attributes (*size*, *number*, and *location* with the *class name*; **P4**, **P5**, and **P6**; see Figure 2).

However, the non-image-specific attributes together with *class name* degrades this performance (**P7**, **P8**, **P9**). Contrastingly for DFU 2022 dataset, prompts **P7**, **P8**, and **P9** achieve the highest performance. This could be due to increased familiarity of pretrained models with general descriptions of feet and skin, available in natural images, compared to colon endoscopy descriptions. Conversely, using just the *class name* (**P1**) or adding a non-image-specific general description (**P7**) seems to perform better than using multiple image-specific attributes in CLIPSeg. This shows that pretraining data and VLSM architecture have a complex relationship with the target medical segmentation task.

**Making prompts richer does not always help during finetuning.** Figure 2 shows that the DSC variation across prompt type is minimal in the finetuned setting for all the models. Prompt with only class name (**P1**) improves segmentation performance in radiology datasets for all four VLSMs. While CRIS' performance almost saturates after adding the class name and mask shape (**P2**), the rest of the models have similar performance for all the prompts except **P0** with multi-class segmentation (CAMUS and CheXlocalize).

BiomedCLIPSeg and BiomedCLIPSeg-D, despite being based on a VLM pretrained on medical data, consistently perform poorly across all prompts compared to CLIP and CLIPSeg. This is likely because it has not been further pretrained for segmentation tasks on a large-scale dataset. As we dig deeper into understanding the impact of individual attributes and robustness of VLSMs, we present the performance of only CLIPSeg and CRIS[5].

**In the finetuned setting, CRIS captures some language semantics well, but CLIPSeg does not.** To check how well the VLSMs captured the semantics of the prompt, we replace attribute values in input prompts during inference with: (**i**) a random and uncommon English word, and (**ii**) semantically wrong or opposite value sampled from the set of possible values of the same attribute (e.g., "small" by "large" for the *size* attribute). The former aims to see how strongly the semantics of words familiar to the model drive output, and the latter allows witnessing the importance of the attribute's presence vs. unfamiliar words' presence. Figure 3 shows results on five datasets with five altered attributes, where CLIPSeg has almost no change in results. To further confirm that there is almost no role of different prompts in CLIPSeg, we found that when testing the CLIPSeg model trained on prompt **P6** by sending only the *class name* (**P1**), the performance was almost the same compared to using **P6** for all radiology datasets.
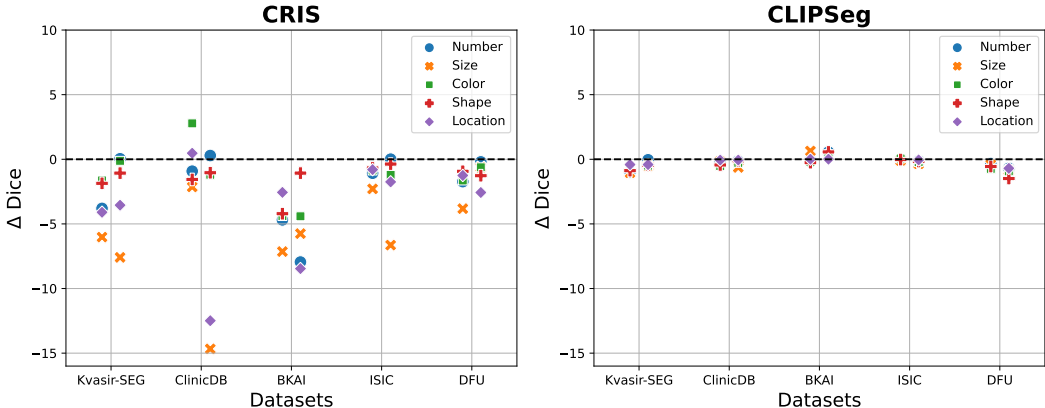


Figure 3: Relative change in DSC when replacing attribute values by a random uncommon English word (left of vertical lines) or semantically wrong opposite value such as replacing 'large' with 'small' (right of vertical lines).

In Figure 3, CRIS's performance drops considerably, with the most significant drop for *size* and *location*. The decline is much more pronounced when giving semantically opposite values compared to random uncommon English words, suggesting that it learned the semantics very well. This is further verified when we look qualitatively into the predicted segmentation masks of CRIS with

---

[5]Additionaly, we have also trained both the models, keeping their encoders frozen whose results are shown in Appendix F.1 of supplementary section

correct vs. incorrect prompts. Figure 4 shows examples of the images having the highest drops in DSC for two datasets when replacing values for the most sensitive attributes [6]



(a) Kvasir-SEG for attribute *size*
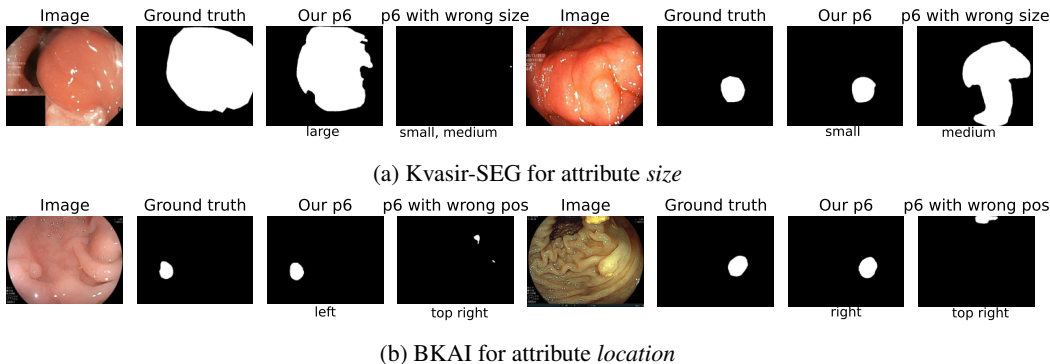


(b) BKAI for attribute *location*

Figure 4: Effect of randomly swapping attribute values on **P6** prompt for non-radiology datasets with another value within the value set of the attributes in the dataset.

**Finetuned VLSMs have comparable performance to SOTA segmentation models**  Table 3 shows results of the experiments carried out to compare VLSMs vs. traditional CNN-based models on their ability to learn in two scenarios: when trained on (**i**) individual specialized datasets or (**ii**) a pooled dataset that combines diverse datasets into a single training set. While the segmentation models (CNNs and VLSMs) achieve better on pooled endoscopy datasets than individual endoscopy datasets, performance mostly drops when training on a pooled set comprising all the datasets. VLSMs outperform image-only off-the-shelf CNN-based methods in most cases compared to our common setting. However, the state-of-the-art results[7] are better, although VLSMs seem to have competitive performance.

Table 3: Performance of VLSMs and CNN models when finetuning in different combinations of datasets. **Bold** for each column shows the best result among all models for the specific dataset combination. **Bold and underline** shows the best result among all the models finetuned for all combinations *i.e., except the state-of-the-art (SOTA)*.

| Tested Dataset → | | Kvasir-SEG | ClinicDB | BKAI | CVC-300 | CVC-ColonDB | ETIS | ISIC | DFU | CAMUS | BUSI | CheXlocalize |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Finetuned Dataset** | **Model** | | | | | | | | | | | |
| **Individual** | CRIS | 86.97 | 84.62 | 84.06 | - | - | - | 89.38 | 71.01 | 88.07 | **83.42** | 57.99 |
| | CLIPSeg | **88.75** | **89.97** | **85.52** | - | - | - | **92.23** | **73.48** | 89.03 | 82.72 | **59.45** |
| | UNet | 84.77 | 85.65 | 83.79 | - | - | - | 90.40 | 67.87 | 90.19 | 75.21 | 50.29 |
| | UNet++ | 84.70 | 84.16 | 84.61 | - | - | - | 90.12 | 69.95 | 89.95 | 72.55 | 49.53 |
| | DeepLabv3+ | 84.11 | 89.11 | 84.95 | - | - | - | 90.66 | 67.89 | **90.43** | 70.57 | 49.95 |
| | *SOTA** | 95.02 | 95.73 | 89.30 | - | - | - | 92.00 | 72.87 | 94.10 | 89.80 | - |
| **Pooled** | CRIS | 86.94 | 84.49 | **82.8** | 82.97 | 68.17 | 60.95 | 91.11 | 67.36 | 81.26 | **78.39** | 56.19 |
| | CLIPSeg | 86.28 | **85.99** | 82.59 | 86.08 | 70.95 | **67.35** | 91.33 | 70.59 | 87.56 | 74.55 | **57.15** |
| | UNet | 36.60 | 26.10 | 37.70 | 4.94 | 8.55 | 12.00 | 64.90 | 38.60 | 76.82 | 44.60 | 38.00 |
| | UNet++ | 80.52 | 78.21 | 77.87 | **87.80** | 51.92 | 48.16 | 88.41 | 65.78 | 89.99 | 75.59 | 53.88 |
| | DeepLabv3+ | 82.40 | 82.70 | 77.60 | 84.40 | 59.30 | 48.30 | 89.60 | 67.70 | **90.17** | 77.80 | 54.56 |
| **Endoscopy Pooled** | CRIS | 89.41 | 85.99 | 85.29 | 86.03 | **74.23** | **71.73** | - | - | - | - | - |
| | CLIPSeg | **89.78** | **90.37** | **86.45** | 88.35 | 73.37 | 70.82 | - | - | - | - | - |
| | UNet | 85.45 | 88.17 | 84.70 | **90.27** | 67.87 | 61.84 | - | - | - | - | - |
| | UNet++ | 83.99 | 85.44 | 82.27 | 89.4 | 66.61 | 55.62 | - | - | - | - | - |
| | DeepLabv3+ | 87.87 | 87.60 | 84.38 | 87.54 | 69.95 | 65.24 | - | - | - | - | - |
| | *SOTA Sources | [11] | [13] | [44] | - | - | - | [18] | [34] | [35] | [58] | - |

**VLSMs adapt better to distribution shifts.**  To assess the ability of the segmentation models to transfer the knowledge learned from one dataset to another similar one, we train the models on each large endoscopy dataset (Kvasir-SEG, ClinicDB, and BKAI) and evaluate them on all the endoscopy

---

[6]More examples provided in the supplementary in Appendix E.

[7]Except for CAMUS and ISIC, may have different training, validation, and test splits due to the unavailability of the standard splits in literature.

datasets. As shown in Table 4, we can observe that VLSMs perform better in all the cases than the conventional models for endoscopic datasets. The performance drops of VLSMs are smaller than that of the conventional models when trained in a different distribution than that of the test set.

Table 4: Segmentation models performance on out-of-distribution endoscopy datasets. **Bold** shows the best result across the model concerning the tested dataset for each finetuning dataset. **<u>Bold and underline</u>** shows the best result for the dataset across the finetuning datasets. The shaded results correspond to results in test sets of the same distribution, while the rest are on out-of-distribution test sets.

| Tested on → Finetuned on ↓ | Model ↓ | Kvasir-SEG | ClinicDB | BKAI | CVC-300 | CVC-ColonDB | ETIS |
|---|---|---|---|---|---|---|---|
| **Kvasir-SEG** | CRIS | 86.97 | 75.47 | 69.46 | 82.47 | 67.51 | 47.84 |
| | CLIPSeg | **<u>88.75</u>** | **80.15** | **77.50** | **85.52** | **<u>71.81</u>** | **<u>65.27</u>** |
| | UNet | 84.77 | 64.84 | 66.22 | 77.16 | 50.81 | 34.98 |
| | UNet++ | 84.70 | 68.15 | 61.76 | 79.35 | 52.3 | 32.81 |
| | DeepLabv3+ | 84.11 | 68.0 | 63.57 | 76.93 | 58.41 | 33.81 |
| **ClinicDB** | CRIS | 79.94 | 84.62 | 69.69 | 84.10 | 67.78 | 53.36 |
| | CLIPSeg | **85.29** | **<u>89.97</u>** | **73.51** | **<u>88.01</u>** | **69.72** | **59.96** |
| | UNet | 65.80 | 85.65 | 35.26 | 73.91 | 55.01 | 29.66 |
| | UNet++ | 61.93 | 84.16 | 38.81 | 71.15 | 55.05 | 23.16 |
| | DeepLabv3+ | 66.63 | 89.11 | 40.89 | 82.05 | 61.79 | 39.53 |
| **BKAI** | CRIS | 79.64 | 72.31 | 84.06 | 85.07 | **64.58** | 61.24 |
| | CLIPSeg | **84.66** | **76.10** | **<u>85.52</u>** | **85.19** | 63.77 | **63.69** |
| | UNet | 68.42 | 62.20 | 83.79 | 60.13 | 44.52 | 42.91 |
| | UNet++ | 70.64 | 62.66 | 84.61 | 82.44 | 55.60 | 46.84 |
| | DeepLabv3+ | 69.02 | 61.99 | 84.95 | 77.47 | 53.15 | 49.61 |

## 5 Discussion

The zero-shot prediction on medical images by VLSMs pretrained on natural images does not provide satisfactory accuracy for any practical use. However, it does provide functional joint text-image representation, which can be further finetuned on triplets of medical images, text, and masks. Our study provides valuable insights into the prompt design, role of various attributes, and models' performance when finetuning for a wide range of datasets. While the best-performing prompts vary with datasets, the best performances are primarily in attributes the models are likely familiar with during pretraining in the natural domain. For example, *size*, *location*, and *number* are well captured by CRIS trained on RefCOCO Dataset [28] for referring image segmentation.

The optimum translation of a pretrained model's ability to leverage such attributes from zero-shot to finetuning settings depends upon multiple factors – the quality and consistency in prompt attributes, the diversity in images concerning the attribute, the knowledge captured in pretrained encoders, and very importantly, the saliency in images corresponding to the semantics. Compared to CRIS, CLIPSeg seemingly fails to capture fine-grained textual information injected by making our prompts richer, indicating that the image encoder's representations are predominantly utilized during finetuning. This could be due to the following two reasons.

1. The notable architectural difference between the two – while CRIS focuses on preserving token-level intervention at refined and granular levels, CLIPSeg uses sentence-level embedding as the intervention from the text encoder.

2. CRIS was trained by updating CLIP's encoders and the segmentation decoder, while CLIPSeg was trained by keeping CLIP frozen. As a result, CRIS's pretrained encoders are familiar with the interaction between image and text pairs tailored for segmentation tasks.

BiomedCLIP, trained on medical test-image pairs, is more familiar with the medical domain than any other encoder or segmentation model encoders we used. In our experiments, though, VLSMs pretrained on natural images performed better than the VLSMs we introduced by adding a decoder to BiomedCLIP. CRIS and CLIPSeg's familiarity with segmentation tasks, albeit in the natural domain, was more valuable than the domain knowledge in BiomedCLIP.

# 6 Limitations

A limited number of VLSMs are trained on large-scale image segmentation data with language prompts. When adapting it to our case, some could not be covered due to a lack of source code, pretrained weights, or reproducibility issues. For example, ZegCLIP gave constant zero scores in zero-shot settings, was trained without background class, and had many channels not amenable to our setup, and SAM does not support text prompts[8]. Nevertheless, as discussed in Section 3, the four VLSMs we considered cover interesting diversity – architectural variation in leveraging global level and token level information in prompts, trained end-to-end for referring image segmentation vs. finetuned only decoder during with segmentation data, based on VLM pretrained on natural vs. medical domain, etc.

Quick glimpses through the output of the VQA model suggest that they are not always reliable. However, we could not manually conduct a thorough quality check to analyze the impact of good vs. bad prompts generated by VQA. Enhancing VQA and Masked Language Models to generate highly reliable automated prompts from specific medical images could enable a more scalable analysis of the kind presented in Figure 3 at a large scale.

Our study focused on 2D medical images and did not include common 3D imaging modalities like MRI or CT scans. While the results seen in ultrasound and X-ray are likely to extend to 2D slices extracted from MRI or CT scans, it is not clear how the results or models extend to 3D volumes. Investigating the adaptation of VLSMs to 3D medical images could open up exciting avenues for future research.

# 7 Perspectives and Conclusions

Although VLSMs utilizing language prompts during segmentation have a big potential for a robust segmentation system, major works remain before the ML community can realize this potential. CRIS's ability to utilize semantics in a finetuned model compared to CLIPSeg's dominant reliance on image features shows that careful designs to learn image-text representation focusing on pixel-token alignment jointly can provide better segmentation models. BiomedCLIPSeg's underperformance, despite having the VLM pretrained on large-scale medical image-text pairs, shows that it is better to further pretrain end-to-end encoder-decoder VLSMs on large-scale natural image-text pairs rather than training medical VLM-based VLSMs on moderately sized datasets. One important line of research seems to be finding ways to generate (potentially synthetic but realistic) large-scale medical image-mask-text triplets and find efficient ways to teach the network concepts represented by language useful for identifying target structure semantics in medical images. A solid foundation VLSM well-versed with concepts specific to medical images and classes, learning joint representation across different problems and datasets, could address numerous challenges in applying deep learning to clinical applications.

Our work is an essential first step in this direction and provides a valuable benchmarking framework, datasets augmented with prompts, and fascinating insights for future investigation.

# References

[1] W. Al-Dhabyani, M. Gomaa, H. Khaled, and A. Fahmy. Dataset of breast ultrasound images. *Data in brief*, 28:104863, 2020. (see pages: 4, 22)

---

[8]Though the paper mentions that text prompts can be added to SAM, its open-sourced implementation does not support text prompts, and there are no pretrained models for text prompts in SAM at the time of this study. Refer to this GitHub issue: `https://github.com/facebookresearch/segment-anything/issues/93`

[2] J. Amin, M. Sharif, M. Yasmin, T. Saba, M. A. Anjum, and S. L. Fernandes. A new approach for brain tumor segmentation and classification based on score level fusion using transfer learning. *Journal of medical systems*, 43:1–16, 2019. (see page: 2)

[3] N. S. An, P. N. Lan, D. V. Hang, D. V. Long, T. Q. Trung, N. T. Thuy, and D. V. Sang. Blazeneo: Blazing fast polyp segmentation and neoplasm detection. *IEEE Access*, 10:43669–43684, 2022. (see pages: 4, 20)

[4] J. Bernal, F. J. Sánchez, G. Fernández-Esparrach, D. Gil, C. Rodríguez, and F. 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. (see pages: 4, 20)

[5] J. Chen, Y. Lu, Q. Yu, X. Luo, E. Adeli, Y. Wang, L. Lu, A. L. Yuille, and Y. Zhou. Transunet: Transformers make strong encoders for medical image segmentation. *arXiv preprint arXiv:2102.04306*, 2021. (see page: 1)

[6] L.-C. Chen, Y. Zhu, G. Papandreou, F. Schroff, and H. Adam. Encoder-decoder with atrous separable convolution for semantic image segmentation. In *Proceedings of the European conference on computer vision (ECCV)*, pages 801–818, 2018. (see page: 15)

[7] L.-C. Chen, P.-C. Kuo, R. Wang, J. Gichoya, and L. A. Celi. Chest x-ray segmentation images based on mimic-cxr. 2022. (see page: 16)

[8] V. Cheplygina, M. de Bruijne, and J. P. Pluim. Not-so-supervised: a survey of semi-supervised, multi-instance, and transfer learning in medical image analysis. *Medical image analysis*, 54:280–296, 2019. (see page: 2)

[9] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 248–255. IEEE, 2009. (see page: 2)

[10] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2020. (see page: 1)

[11] R.-G. Dumitru, D. Peteleaza, and C. Craciun. Using duck-net for polyp image segmentation. *Scientific Reports*, 13(1):9803, 2023. (see page: 8)

[12] S. Eslami, C. Meinel, and G. De Melo. Pubmedclip: How much does clip benefit visual question answering in the medical domain? In *Findings of the Association for Computational Linguistics: EACL 2023*, pages 1151–1163, 2023. (see page: 3)

[13] K. Fitzgerald and B. Matuszewski. Fcb-swinv2 transformer for polyp segmentation. *arXiv preprint arXiv:2302.01027*, 2023. (see page: 8)

[14] A. Fürst, E. Rumetshofer, J. Lehner, V. T. Tran, F. Tang, H. Ramsauer, D. Kreil, M. Kopp, G. Klambauer, A. Bitto, et al. Cloob: Modern hopfield networks with infoloob outperform clip. *Advances in neural information processing systems*, 35:20450–20468, 2022. (see page: 2)

[15] M. Ghafoorian, A. Mehrtash, T. Kapur, N. Karssemeijer, E. Marchiori, M. Pesteie, C. R. Guttmann, F.-E. de Leeuw, C. M. Tempany, B. Van Ginneken, et al. Transfer learning for domain adaptation in mri: Application in brain lesion segmentation. In *Medical Image Computing and Computer Assisted Intervention- MICCAI 2017: 20th International Conference, Quebec City, QC, Canada, September 11-13, 2017, Proceedings, Part III 20*, pages 516–524. Springer, 2017. (see page: 2)

[16] Y. Gu, R. Tinn, H. Cheng, M. Lucas, N. Usuyama, X. Liu, T. Naumann, J. Gao, and H. Poon. Domain-specific language model pretraining for biomedical natural language processing. *ACM Transactions on Computing for Healthcare (HEALTH)*, 3(1):1–23, 2021. (see page: 5)

[17] D. Gutman, N. C. Codella, E. Celebi, B. Helba, M. Marchetti, N. Mishra, and A. Halpern. 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. (see pages: 4, 20)

[18] M. K. Hasan, M. T. E. Elahi, M. A. Alam, M. T. Jawad, and R. Martí. Dermoexpert: Skin lesion classification using a hybrid convolutional neural network through segmentation, transfer learning, and augmentation. *Informatics in Medicine Unlocked*, 28:100819, 2022. (see page: 8)

[19] A. Hatamizadeh, Y. Tang, V. Nath, D. Yang, A. Myronenko, B. Landman, H. R. Roth, and D. Xu. Unetr: Transformers for 3d medical image segmentation. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, pages 574–584, 2022. (see page: 1)

[20] M. Havaei, A. Davy, D. Warde-Farley, A. Biard, A. Courville, Y. Bengio, C. Pal, P.-M. Jodoin, and H. Larochelle. Brain tumor segmentation with deep neural networks. *Medical image analysis*, 35:18–31, 2017. (see page: 1)

[21] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. (see page: 15)

[22] Z. Huang, Z. Zeng, B. Liu, D. Fu, and J. Fu. Pixel-bert: Aligning image pixels with text by deep multi-modal transformers. *arXiv preprint arXiv:2004.00849*, 2020. (see page: 2)

[23] F. Isensee, P. F. Jaeger, S. A. Kohl, J. Petersen, and K. H. Maier-Hein. nnu-net: a self-configuring method for deep learning-based biomedical image segmentation. *Nature methods*, 18(2):203–211, 2021. (see page: 1)

[24] D. Jha, P. H. Smedsrud, M. A. Riegler, P. Halvorsen, T. de Lange, D. Johansen, and H. 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*, pages 451–462. Springer, 2020. (see pages: 4, 20)

[25] C. Jia, Y. Yang, Y. Xia, Y.-T. Chen, Z. Parekh, H. Pham, Q. Le, Y.-H. Sung, Z. Li, and T. Duerig. Scaling up visual and vision-language representation learning with noisy text supervision. In *International Conference on Machine Learning*, pages 4904–4916. PMLR, 2021. (see pages: 2, 3)

[26] A. E. Johnson, T. J. Pollard, S. J. Berkowitz, N. R. Greenbaum, M. P. Lungren, C.-y. Deng, R. G. Mark, and S. Horng. MIMIC-CXR, a de-identified publicly available database of chest radiographs with free-text reports. *Scientific data*, 6(1):317, 2019. (see pages: 16, 19)

[27] A. E. Johnson, T. J. Pollard, N. R. Greenbaum, M. P. Lungren, C.-y. Deng, Y. Peng, Z. Lu, R. G. Mark, S. J. Berkowitz, and S. Horng. MIMIC-CXR-JPG, a large publicly available database of labeled chest radiographs. *arXiv preprint arXiv:1901.07042*, 2019. (see page: 16)

[28] S. Kazemzadeh, V. Ordonez, M. Matten, and T. Berg. Referitgame: Referring to objects in photographs of natural scenes. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pages 787–798, 2014. (see pages: 4, 9)

[29] C. Kendrick, B. Cassidy, J. M. Pappachan, C. O'Shea, C. J. Fernandez, E. Chacko, K. Jacob, N. D. Reeves, and M. H. Yap. Translating clinical delineation of diabetic foot ulcers into machine interpretable segmentation. *arXiv preprint arXiv:2204.11618*, 2022. (see pages: 4, 20)

[30] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. (see page: 5)

[31] S. Leclerc, E. Smistad, J. Pedrosa, A. Østvik, F. Cervenansky, F. Espinosa, T. Espeland, E. A. R. Berg, P.-M. Jodoin, T. Grenier, et al. Deep learning for segmentation using an open large-scale dataset in 2d echocardiography. *IEEE transactions on medical imaging*, 38(9):2198–2210, 2019. (see pages: 4, 21)

[32] J. Li, D. Li, C. Xiong, and S. Hoi. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *International Conference on Machine Learning*, pages 12888–12900. PMLR, 2022. (see page: 3)

[33] Y. Li, F. Liang, L. Zhao, Y. Cui, W. Ouyang, J. Shao, F. Yu, and J. Yan. Supervision exists everywhere: A data efficient contrastive language-image pre-training paradigm. In *International Conference on Learning Representations*, 2021. (see pages: 2, 3)

[34] T.-Y. Liao, C.-H. Yang, Y.-W. Lo, K.-Y. Lai, P.-H. Shen, and Y.-L. Lin. Hardnet-dfus: Enhancing backbone and decoder of hardnet-mseg for diabetic foot ulcer image segmentation. In *Diabetic Foot Ulcers Grand Challenge*, pages 21–30. Springer, 2022. (see page: 8)

[35] H. J. Ling, D. Garcia, and O. Bernard. Reaching intra-observer variability in 2-d echocardiographic image segmentation with a simple u-net architecture. In *IEEE International Ultrasonics Symposium (IUS)*, 2022. (see page: 8)

[36] T. Lüddecke and A. Ecker. Image segmentation using text and image prompts. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7086–7096, 2022. (see pages: 2, 3)

[37] F. Milletari, N. Navab, and S.-A. Ahmadi. V-net: Fully convolutional neural networks for volumetric medical image segmentation. In *2016 fourth international conference on 3D vision (3DV)*, pages 565–571. Ieee, 2016. (see page: 1)

[38] P. Ngoc Lan, N. S. An, D. V. Hang, D. V. Long, T. Q. Trung, N. T. Thuy, and D. V. Sang. Neounet: Towards accurate colon polyp segmentation and neoplasm detection. In *Advances in Visual Computing: 16th International Symposium, ISVC 2021, Virtual Event, October 4-6, 2021, Proceedings, Part II*, pages 15–28. Springer, 2021. (see pages: 4, 20)

[39] O. Oktay, J. Schlemper, L. Le Folgoc, M. Lee, M. Heinrich, K. Misawa, K. Mori, S. McDonagh, N. Y. Hammerla, B. Kainz, et al. Attention u-net: Learning where to look for the pancreas. In *Medical Imaging with Deep Learning*, 2022. (see page: 1)

[40] Z. Qin, H. H. Yi, Q. Lao, and K. Li. Medical image understanding with pretrained vision language models: A comprehensive study. In *The Eleventh International Conference on Learning Representations*, 2022. (see pages: 2, 5)

[41] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. (see pages: 2, 3)

[42] Y. Rao, W. Zhao, G. Chen, Y. Tang, Z. Zhu, G. Huang, J. Zhou, and J. Lu. Denseclip: Language-guided dense prediction with context-aware prompting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18082–18091, 2022. (see pages: 2, 3)

[43] O. Ronneberger, P. Fischer, and T. Brox. U-net: Convolutional networks for biomedical image segmentation. In *Medical Image Computing and Computer-Assisted Intervention–MICCAI 2015: 18th International Conference, Munich, Germany, October 5-9, 2015, Proceedings, Part III 18*, pages 234–241. Springer, 2015. (see pages: 1, 15)

[44] D. Sang. BKAI-IGH NeoPolyp, 2021. URL `https://kaggle.com/competitions/bkai-igh-neopolyp`. (see page: 8)

[45] A. Saporta, X. Gui, A. Agrawal, A. Pareek, S. Q. Truong, C. D. Nguyen, V.-D. Ngo, J. Seekins, F. G. Blankenberg, A. Y. Ng, et al. Benchmarking saliency methods for chest x-ray interpretation. *Nature Machine Intelligence*, 4(10):867–878, 2022. (see pages: 4, 21)

[46] C. Seibold, S. Reiß, M. S. Sarfraz, R. Stiefelhagen, and J. Kleesiek. Breaking with fixed set pathology recognition through report-guided contrastive training. In *Medical Image Computing and Computer Assisted Intervention–MICCAI 2022: 25th International Conference, Singapore, September 18–22, 2022, Proceedings, Part V*, pages 690–700. Springer, 2022. (see page: 3)

[47] J. Silva, A. Histace, O. Romain, X. Dray, and B. Granado. Toward embedded detection of polyps in wce images for early diagnosis of colorectal cancer. *International journal of computer assisted radiology and surgery*, 9:283–293, 2014. (see pages: 4, 20)

[48] A. Singh, R. Hu, V. Goswami, G. Couairon, W. Galuba, M. Rohrbach, and D. Kiela. Flava: A foundational language and vision alignment model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15638–15650, 2022. (see pages: 2, 3)

[49] N. Tajbakhsh, S. R. Gurudu, and J. Liang. Automated polyp detection in colonoscopy videos using shape and context information. *IEEE transactions on medical imaging*, 35(2):630–644, 2015. (see pages: 4, 20)

[50] N. K. Tomar, D. Jha, U. Bagci, and S. Ali. Tganet: text-guided attention for improved polyp segmentation. In *Medical Image Computing and Computer Assisted Intervention–MICCAI 2022: 25th International Conference, Singapore, September 18–22, 2022, Proceedings, Part III*, pages 151–160. Springer, 2022. (see pages: 3, 5)

[51] D. Vázquez, J. Bernal, F. J. Sánchez, G. Fernández-Esparrach, A. M. López, A. Romero, M. Drozdzal, and A. Courville. A benchmark for endoluminal scene segmentation of colonoscopy images. *Journal of healthcare engineering*, 2017, 2017. (see pages: 4, 20)

[52] Z. Wang, Y. Lu, Q. Li, X. Tao, Y. Guo, M. Gong, and T. Liu. Cris: Clip-driven referring image segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11686–11695, 2022. (see pages: 2, 3)

[53] Z. Wang, Z. Wu, D. Agarwal, and J. Sun. Medclip: Contrastive learning from unpaired medical images and text. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 3876–3887, 2022. (see page: 3)

[54] S. Wazir and M. M. Fraz. Histoseg: Quick attention with multi-loss function for multi-structure segmentation in digital histology images. In *2022 12th International Conference on Pattern Recognition Systems (ICPRS)*, pages 1–7. IEEE, 2022. (see page: 1)

[55] C. Wu, Z. Lin, S. Cohen, T. Bui, and S. Maji. Phrasecut: Language-based image segmentation in the wild. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10216–10225, 2020. (see page: 4)

[56] C. Wu, X. Zhang, Y. Zhang, Y. Wang, and W. Xie. Medklip: Medical knowledge enhanced language-image pre-training. *medRxiv*, pages 2023–01, 2023. (see page: 3)

[57] X. Zhai, X. Wang, B. Mustafa, A. Steiner, D. Keysers, A. Kolesnikov, and L. Beyer. Lit: Zero-shot transfer with locked-image text tuning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18123–18133, 2022. (see pages: 2, 3)

[58] S. Zhang, M. Liao, J. Wang, Y. Zhu, Y. Zhang, J. Zhang, R. Zheng, L. Lv, D. Zhu, H. Chen, et al. Fully automatic tumor segmentation of breast ultrasound images with deep learning. *Journal of Applied Clinical Medical Physics*, 24(1):e13863, 2023. (see page: 8)

[59] S. Zhang, Y. Xu, N. Usuyama, J. Bagga, R. Tinn, S. Preston, R. Rao, M. Wei, N. Valluri, C. Wong, et al. Large-scale domain-specific pretraining for biomedical vision-language processing. *arXiv preprint arXiv:2303.00915*, 2023. (see pages: 2, 3, and 4)

[60] Y. Zhang, H. Jiang, Y. Miura, C. D. Manning, and C. P. Langlotz. Contrastive learning of medical visual representations from paired images and text. In *Machine Learning for Healthcare Conference*, pages 2–25. PMLR, 2022. (see page: 3)

[61] Z. Zhou, M. M. Rahman Siddiquee, N. Tajbakhsh, and J. Liang. Unet++: A nested u-net architecture for medical image segmentation. In *Deep Learning in Medical Image Analysis and Multimodal Learning for Clinical Decision Support: 4th International Workshop, DLMIA 2018, and 8th International Workshop, ML-CDS 2018, Held in Conjunction with MICCAI 2018, Granada, Spain, September 20, 2018, Proceedings 4*, pages 3–11. Springer, 2018. (see pages: 1, 15)

[62] Z. Zhou, Y. Lei, B. Zhang, L. Liu, and Y. Liu. Zegclip: Towards adapting clip for zero-shot semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11175–11185, 2023. (see page: 3)

# A  Impact on Society

The incorporation of language prompts in medical image segmentation has the potential to impact society, particularly in clinical settings, significantly. By enabling radiologists to quickly and accurately segment complex shapes using just a few words, language prompts offer a more interpretable and explainable approach compared to traditional visual prompts such as points or boxes.

One significant advantage of language prompts is their ability to convey detailed information about normal and abnormal structures' texture, shape, and spatial relationships. This allows for a more comprehensive medical image understanding, facilitating more accurate segmentation results. Additionally, language prompts can be easily adapted to new classes, making them highly versatile and adaptable in various medical scenarios.

Using language prompts in medical image segmentation can improve the efficiency and effectiveness of radiologists' work, potentially leading to faster diagnoses and treatment decisions. Moreover, the interpretability of language prompts can aid in building trust and confidence among healthcare professionals and patients as the reasoning behind the segmentation process becomes more transparent.

Overall, the integration of language prompts in medical image segmentation has the potential to revolutionize clinical practices, providing radiologists with a powerful tool to enhance their segmentation capabilities and ultimately improve patient care outcomes.

We strongly encourage and invite other researchers to contribute to this field of study. This research paper has no negative impact on society or further research in medical imaging, as we have adhered to ethical considerations in medical imaging and have not expressed disapproval of any previous studies.

# B  Dataset and Code Access

The GitHub repository[9] contains the source code with detailed documentation, the generated prompts for all the datasets, and thorough instructions along with the relevant links to access the individual image-mask pair datasets used in this work.

# C  Experiments

## C.1  VLSM Finetuning Experiments

For the five non-radiology datasets (Kvasir-SEG, ClinicDB, BKAI, ISIC, and DFU), we finetune VLSMs with ten prompts for an individual dataset, resulting in 50 experiments for each VLSM. Similarly, in the case of radiology datasets (CAMUS, BUSI, and CheXlocalize), we have a total of 22 finetuning experiments for each VLSM. We also finetune CRIS and CLIPSeg with the pooled datasets comprising only endoscopic and all datasets. Thus, including all varieties with the VLSMs and the different prompting mechanisms, we have 442 finetuning experiments.

The average time to fine-tune CRIS for a dataset on a prompt is approximately 60 minutes in our training setup, running 45 epochs on average. For CLIPSeg, the average training time is 40 minutes, running for 90 epochs on average. BiomedCLIPSeg's and BiomedCLIPSeg-D's average training times are 20 minutes and 30 minutes, running for 80 epochs and 50 epochs, respectively. We monitored the segmentation metric on the held-out validation sets for early stopping, with patience of 50 epochs for CLIPSeg variants and 10 epochs for CRIS.

## C.2  CNN-based Experiments

For comparative analysis, we consider three of the conventional CNN-based segmentation models: UNet [43], UNet++[61], and DeepLabV3+ [6]. For all of the models, we use pretrained ResNet-50 [21] as the backbone, and default parameters given by the framework *Segmentation Models Pytorch*

---

[9]https://github.com/naamiinepal/medvlsm

are chosen as the model hyperparameters. We use Dice loss for error propagation within the models with Adam optimizer of learning rate $10^{-3}$ and zero weight decay.

## D    PubMedBERT's failure to give reliable output

Table 5 contains the predictions of PubMedBERT for the masked language modeling in different datasets.

Table 5: PubMedBERT's top five predictions for the masked language modeling inference. The predictions are ordered in the descending order of the probability generated by the model. The model has high uncertainty as the maximum probability is about $0.1$. The predictions are almost the same and uninformative, which is more prominent in the radiology datasets.

| Dataset | Masked sentence | Top-5 Predictions |
|---|---|---|
| All Endoscopy* | The location of the polyp is [MASK]. | [variable, unknown, varied, unpredictable, uncertain] |
| | Polyp is located at [MASK]. | [bifurcation, apex, rectum, midline, right] |
| | The shape of polyp is [MASK]. | [irregular, variable, oval, round, different] |
| | Polyp is [MASK] in shape. | [oval, irregular, round, spherical, cylindrical] |
| | The color of the polyp is [MASK]. | [yellow, red, blue, brown, pink] |
| | Polyp is [MASK] in color. | [yellow, white, red, black, green] |
| ISIC | The location of skin melanoma is [MASK]. | [unknown, variable, unusual, unpredictable, rare] |
| | The color of skin melanoma is [MASK]. | [red, yellow, brown, black, blue] |
| | Skin melanoma is [MASK] in texture. | [heterogeneous, variable, soft, irregular, fibrous] |
| | Skin cancer is located at [MASK]. | [extremities, birth, puberty, adolescence, skin] |
| | Skin cancer is [MASK] in texture. | [heterogeneous, unique, variable, diverse, distinctive] |
| DFU | The location of a diabetic foot ulcer is at [MASK]. | [first, rest, ankle, home, foot] |
| | Diabetic foot ulcer is located at [MASK]. | [ankle, heel, foot, extremities, feet] |
| | The location of the foot ulcer is [MASK]. | [ankle, knee, first, heel, night] |
| | Foot ulcer is located at [MASK]. | [ankle, heel, foot, knee, night] |
| CAMUS | The left ventricular cavity is [MASK] in shape. | [spherical, triangular, normal, oval, round] |
| | The myocardium is [MASK] in shape. | [spherical, cylindrical, circular, round, triangular] |
| | The left atrium cavity is [MASK] in shape. | [oval, round, triangular, spherical, irregular] |
| | The left ventricular cavity is located at [MASK]. | [diastole, apex, rest, $90°$, $45°$] |
| | The myocardium is located at [MASK]. | [rest, apex, risk, diastole, birth] |
| | The left atrium cavity is located at [MASK]. | [diastole, right, left, $90°$, apex] |
| BUSI | The malignant breast tumor is [MASK] in shape. | [round, irregular, oval, solid, spherical] |
| | The benign breast tumor is [MASK] in shape. | [oval, round, irregular, solid, spherical] |
| CheXlocalize | Airspace Opacity is [MASK] in shape. | [irregular, oval, round, triangular, globular] |
| | Enlarged Cardiomediastinum is [MASK] in shape. | [oval, triangular, irregular, round, rounded] |
| | Cardiomegaly is [MASK] in shape. | [irregular, triangular, normal, oval, round] |
| | Lung Opacity is [MASK] in shape. | [irregular, round, oval, nodular, reticular] |
| | Consolidation is [MASK] in shape. | [spherical, circular, triangular, irregular, round] |
| | Atelectasis is [MASK] in shape. | [irregular, oval, triangular, spherical, round] |
| | Pleural Effusion is [MASK] in shape. | [irregular, round, oval, spherical, solid] |

*This includes six datasets of endoscopy: Kvasir-SEG, ClinicDB, BKAI, CVC-300, CVC-ColonDB, ETIS

## E    Some visualizations and qualitative analysis

Some visualizations and qualitative analysis are shown in Figures 5 and 6.

## F    Results

### F.1    Finetuning only the Decoders for CLIP-based VLSMs

Tables 6 and 7 show the results of VLSMs with finetuned the decoder while keeping the encoders frozen.

### F.2    Using radiology reports for lung segmentation

To examine the usage of free-text radiology reports of chest x-rays for segmentation, we utilize 1,141 frontal-view CXRs randomly selected from the MIMIC-CXR database [7, 26, 27]. This dataset contains the segmentation of lungs, which has been verified manually. We use the free-text radiology

---

[10]https://github.com/qubvel/segmentation_models.pytorch

(a) Kvasir-SEG for attribute *size*

(b) Bkai for attribute *location*

(c) ClinicDB for attribute *size*

(d) DFU for attribute *location*
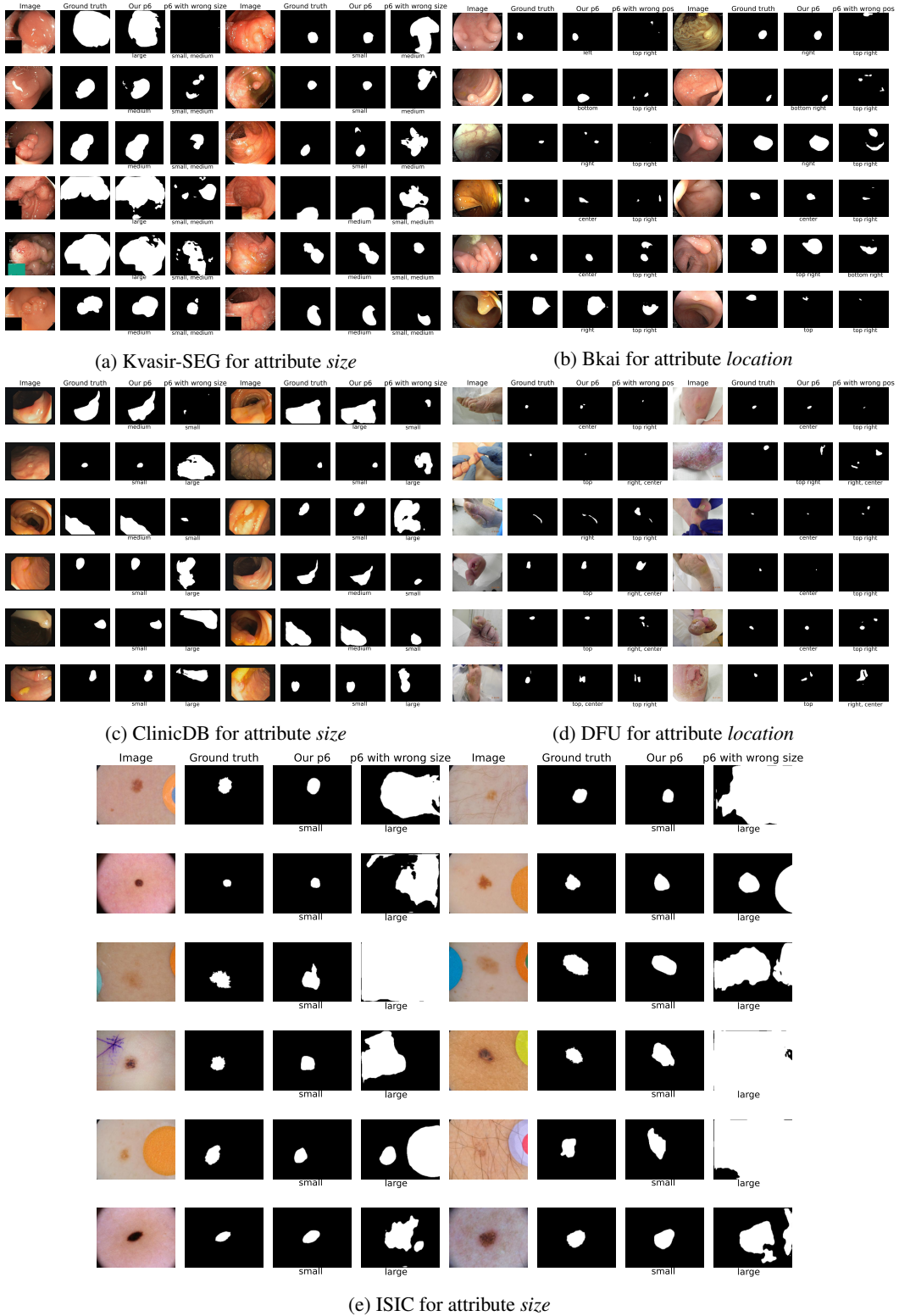
(e) ISIC for attribute *size*

Figure 5: Visualization of CRIS's performance when prompt attributes are changed using a wrong attribute value. For each medical image, three corresponding masks are displayed: ground truth mask, output mask for the corresponding prompt, and output mask after altering an attribute value of the prompts.

| Inputs | Ground Truth | CRIS Prediction | CLIPSeg Prediction |
|---|---|---|---|

one small pink round
polyp, located in bottom
right of the image

medium brown circular
skin melanoma

one medium pink oval
foot ulcer, located in
top of the image

Myocardium in four-cham-
ber view of the heart
at end of the diastole
cycle of a male.

One medium irregular-
shaped tumor in the
breast ultrasound image.

Enlarged Cardiomedias-
tinum of shape rectangle,
and located in center of the
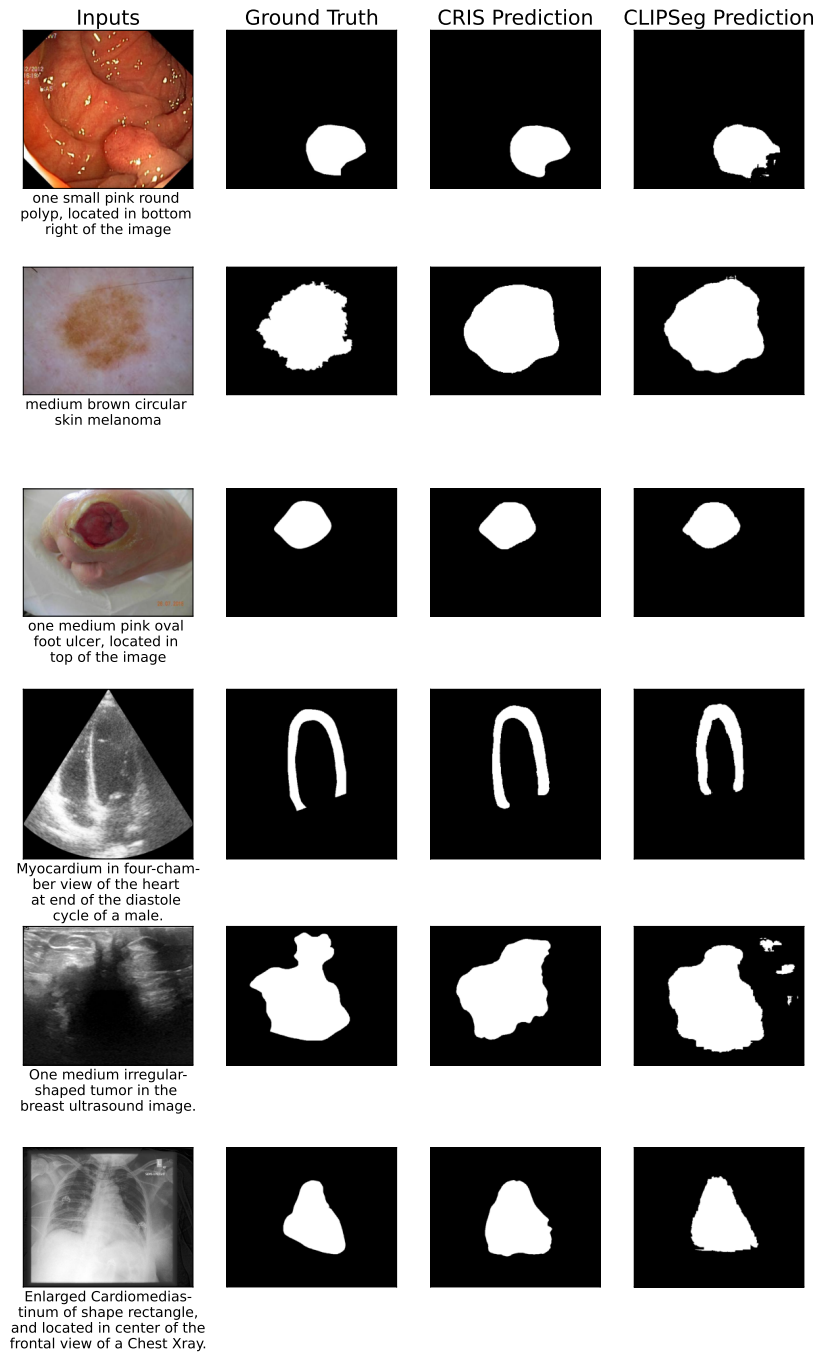frontal view of a Chest Xray.

Figure 6: Sample input, ground truth, and models' predictions

Table 6: Finetuned segmentation dice score of CRIS on different datasets on different sets of prompts with frozen CLIP.

| Prompt → / Dataset ↓ | P0 | P1 | P2 | P3 | P4 | P5 | P6 | P7 | P8 | P9 |
|---|---|---|---|---|---|---|---|---|---|---|
| **Kvasir-SEG** | $52.72_{\pm24.42}$ | $60.15_{\pm27.27}$ | $82.35_{\pm20.55}$ | $82.47_{\pm21.33}$ | $84.81_{\pm18.89}$ | $85.73_{\pm17.05}$ | $\mathbf{86.93}_{\pm15.28}$ | $84.94_{\pm18.08}$ | $86.51_{\pm15.87}$ | $86.62_{\pm15.54}$ |
| **ClinicDB** | $50.84_{\pm34.25}$ | $52.92_{\pm32.93}$ | $80.61_{\pm24.48}$ | $81.62_{\pm23.48}$ | $84.37_{\pm20.02}$ | $\mathbf{84.54}_{\pm19.56}$ | $81.41_{\pm22.38}$ | $83.25_{\pm23.44}$ | $83.66_{\pm22.67}$ | $83.68_{\pm22.31}$ |
| **BKAI** | $50.99_{\pm34.44}$ | $74.88_{\pm30.69}$ | $82.24_{\pm23.31}$ | $83.24_{\pm24.39}$ | $82.74_{\pm26.08}$ | $84.01_{\pm23.59}$ | $\mathbf{85.38}_{\pm21.28}$ | $81.22_{\pm25.14}$ | $83.35_{\pm23.92}$ | $85.13_{\pm22.24}$ |
| **ISIC** | $6.38_{\pm12.17}$ | $0.78_{\pm3.75}$ | $88.98_{\pm13.37}$ | $89.13_{\pm12.67}$ | $90.30_{\pm11.35}$ | $\mathbf{90.58}_{\pm10.45}$ | $90.11_{\pm11.76}$ | $89.89_{\pm11.88}$ | $90.42_{\pm10.62}$ | $90.14_{\pm11.41}$ |
| **DFU** | $25.29_{\pm23.53}$ | $0.00_{\pm0.00}$ | $68.07_{\pm29.43}$ | $66.84_{\pm30.04}$ | $68.87_{\pm27.7}$ | $68.39_{\pm28.9}$ | $67.12_{\pm29.98}$ | $\mathbf{69.00}_{\pm28.89}$ | $68.74_{\pm28.81}$ | $67.53_{\pm30.61}$ |
| **CAMUS** | $8.29_{\pm11.90}$ | $87.01_{\pm10.24}$ | $87.04_{\pm10.54}$ | $86.82_{\pm11.87}$ | $86.77_{\pm12.18}$ | $86.71_{\pm11.99}$ | $86.71_{\pm12.09}$ | $\mathbf{87.26}_{\pm10.56}$ | N/A | N/A |
| **BUSI** | $22.41_{\pm25.73}$ | $76.98_{\pm29.38}$ | $80.18_{\pm25.70}$ | $79.79_{\pm25.95}$ | $79.15_{\pm27.23}$ | $81.49_{\pm23.59}$ | $\mathbf{81.78}_{\pm22.74}$ | N/A | N/A | N/A |
| **CheXlocalize** | $8.94_{\pm12.81}$ | $54.85_{\pm26.71}$ | $54.75_{\pm27.1}$ | $54.6_{\pm26.64}$ | $56.69_{\pm26.16}$ | $56.89_{\pm26.11}$ | $54.22_{\pm26.76}$ | N/A | N/A | N/A |

Table 7: Finetuned segmentation dice score of CLIPSeg on different datasets on different sets of prompts with frozen CLIP.

| Prompt → / Dataset ↓ | P0 | P1 | P2 | P3 | P4 | P5 | P6 | P7 | P8 | P9 |
|---|---|---|---|---|---|---|---|---|---|---|
| **Kvasir-SEG** | $87.93_{\pm14.93}$ | $86.28_{\pm17.65}$ | $86.8_{\pm15.57}$ | $88.02_{\pm13.45}$ | $87.29_{\pm15.71}$ | $88.14_{\pm13.18}$ | $88.69_{\pm12.29}$ | $88.11_{\pm13.24}$ | $\mathbf{89.12}_{\pm12.63}$ | $86.68_{\pm16.82}$ |
| **ClinicDB** | $89.14_{\pm12.32}$ | $88.89_{\pm12.53}$ | $89.43_{\pm9.66}$ | $89.85_{\pm10.87}$ | $89.33_{\pm10.84}$ | $89.13_{\pm11.49}$ | $89.65_{\pm10.8}$ | $88.45_{\pm12.58}$ | $\mathbf{89.98}_{\pm9.71}$ | $87.48_{\pm13.79}$ |
| **BKAI** | $83.27_{\pm18.11}$ | $83.8_{\pm19.36}$ | $83.35_{\pm18.65}$ | $83.82_{\pm20.13}$ | $84.64_{\pm17.29}$ | $85.25_{\pm14.61}$ | $84.96_{\pm16.12}$ | $\mathbf{85.39}_{\pm14.74}$ | $84.65_{\pm16.86}$ | $85.09_{\pm16.43}$ |
| **ISIC** | $91.69_{\pm8.67}$ | $92.03_{\pm7.3}$ | $91.93_{\pm7.78}$ | $91.48_{\pm9.36}$ | $91.98_{\pm7.11}$ | $92.05_{\pm6.92}$ | $92.16_{\pm7.07}$ | $91.64_{\pm8.26}$ | $\mathbf{92.24}_{\pm6.98}$ | $92.1_{\pm8.16}$ |
| **DFU** | $72.78_{\pm25.38}$ | $72.0_{\pm25.43}$ | $72.65_{\pm24.31}$ | $72.72_{\pm24.88}$ | $72.13_{\pm24.81}$ | $72.72_{\pm24.56}$ | $72.6_{\pm24.84}$ | $71.92_{\pm25.42}$ | $73.22_{\pm24.51}$ | $\mathbf{73.58}_{\pm24.44}$ |
| **CAMUS** | $46.51_{\pm13.64}$ | $88.6_{\pm7.46}$ | $88.85_{\pm7.18}$ | $88.59_{\pm7.44}$ | $88.87_{\pm7.03}$ | $88.74_{\pm6.98}$ | $\mathbf{88.93}_{\pm7.19}$ | $88.37_{\pm8.47}$ | N/A | N/A |
| **BUSI** | $74.04_{\pm33.61}$ | $80.46_{\pm26.85}$ | $80.98_{\pm25.58}$ | $80.76_{\pm26.18}$ | $80.53_{\pm28.05}$ | $68.88_{\pm36.09}$ | $\mathbf{82.25}_{\pm24.1}$ | N/A | N/A | N/A |
| **CheXlocalize** | $45.52_{\pm25.11}$ | $58.39_{\pm24.62}$ | $57.42_{\pm25.31}$ | $58.54_{\pm25.01}$ | $\mathbf{58.87}_{\pm25.41}$ | $58.47_{\pm25.41}$ | $58.27_{\pm24.71}$ | N/A | N/A | N/A |

reports provided in the MIMIC-CXR Database [26] as the only prompt (P1), and the results are reported in Table 8.

Table 8: Zero-shot and finetuning dice scores of the CRIS and CLIPSeg Manually labeled Chest X-ray Segmentation Dataset. We have used the actual radiology reports as **P1**. P0 indicates an empty prompt.

| Models ↓ | Experiment ↓     Prompt → | P0 | P1 |
|---|---|---|---|
| **CRIS** | **Zero-shot** | $44.8_{\pm18.97}$ | $40.73_{\pm18.95}$ |
| | **Finetuning** | $81.66_{\pm5.65}$ | $90.99_{\pm1.41}$ |
| **CLIPSeg** | **Zero-shot** | $0.26_{\pm2.35}$ | $0.09_{\pm0.88}$ |
| | **Finetuning** | $91.39_{\pm1.09}$ | $91.22_{\pm1.26}$ |

# G   Prompt Composition

The prompts used during the training for various datasets are shown below. If there are multiple templates for the same prompts for a dataset, one is randomly chosen during the training to increase the regularization for the models.

## G.1   Non-radiology images

### G.1.1   Endoscopy Datasets

A total of six endoscopy datasets (polyp segmentation image-mask pairs) have been used for finetuning and evaluating our proposed models: Kvasir-SEG [24], ClinicDB [4], BKAI [3, 38], CVC-300 [51], CVC-ColonDB [49], and ETIS [47]. The last three datasets have few numbers of image-masks pairs, so they are used only for testing and evaluating the trained models.

1. **P0**: "" (No prompt)
2. **P1**: "*class name*"
     • *polyp*
3. **P2**: "*shape class name*"
     • *round polyp*
4. **P3**: "*color shape class name*"
     • *pink round polyp*
5. **P4**: "*size color shape class name*"
     • *medium pink round polyp*
6. **P5**: "*number size color shape class name*"
     • *one medium pink round polyp*
7. **P6**: "*number size color shape class name*, located in the *location* of the image"
     • *one medium pink round polyp*, located in the *top left* of the image
8. **P7**: "*class name*, which is a *general description of the class*"
     • *polyp*, which is a *small lump in the lining of colon*
9. **P8**: "*number size color shape class name*, which is a *general description of the class*"
     • *one medium pink round polyp*, which is a *small lump in the lining of colon*
10. **P9**: "*number size color shape class name*, which is a *general description of the class* located in the *location* of the image "
     • *one medium pink round polyp*, which is a *small lump in the lining of colon* located in the *top left* of the image

For *General Description of the class*, prompts were built using information about the subject on the internet. Five such descriptions were designed for each dataset, and one random sample was selected each time as the *general description of the class* attribute whenever the prompts **p7**, **p8**, and **p9** were used.

### G.1.2   ISIC and DFU-2022

The templates of prompts for the DFU-2022 [29] and ISIC [17] datasets used were the same as the above examples for endoscopy images, with *class name* and *general description of the class* being different. We used class names **skin melanoma** and **foot ulcer** for the two datasets, respectively.

The five *General Description of the class* for each of the three types of photographic datasets used is listed in the table below.

Table 9: General Descriptions selected for each of the photographic datsaets

| Endoscopy Datasets | ISIC | DFU-2022 |
|---|---|---|
| → a projecting growth of tissue | → a spot with dark speckles | → a wound in foot and toes |
| → often a bumpy flesh in rectum | → a spot with irregular texture | → a sore in foot and toes |
| → a small lump in the lining of colon | → a dark sore with irregular texture | → a sore in skin of foot and toe |
| → a tissue growth that often resemble mushroom-like stalks | → an irregular sore with speckles | → an abnormality in foot and toes |
| → an abnormal growth of tissues projecting from a mucous membrane | → a rough wound on skin | → an open sore or lesion in foot and toes |

## G.2 Radiology Images

### G.2.1 CheXlocalize

The prompts for the CheXlocalize [45] dataset are listed below.

1. **P0**: "" (No prompt)
2. **P1**: "*labels* in a chest Xray."
   - *Airspace Opacity* in a chest Xray.
3. **P2**: "*labels* in the *xray_view* view of a Chest Xray."
   - Airspace Opacity in the *frontal* view of a Chest Xray.
4. **P3**: "*labels* of shape *shape* in the *xray_view* view of a Chest Xray."
   - Airspace Opacity of shape *rectangle* in the frontal view of a Chest Xray.
5. **P4**: "*labels* of shape *shape*, and located in *location* of the *xray_view* view of a Chest Xray."
   - Airspace Opacity of shape rectangle, and located in *right* of the frontal view of a Chest Xray.
6. **P5**: "*labels* of shape *shape*, and located in *location* of the *xray_view* view of a Chest Xray. *pathology* are present."
   - Airspace Opacity of shape rectangle, and located in right of the frontal view of a Chest Xray. *Enlarged Cardiomediastinum, Cardiomegaly, Lung Opacity, Consolidation, Atelectasis, Pleural Effusion* are present.
7. **P6**: "*labels* in a Chest Xray. *pathology* are present."
   - Airspace Opacity in a Chest Xray. Enlarged Cardiomediastinum, Cardiomegaly, Lung Opacity, Consolidation, Atelectasis, Pleural Effusion are present.

### G.2.2 CAMUS

The prompts for the CAMUS [31] dataset are listed below.

1. Class of Current Image
   - *Left ventricular cavity*, *Myocardium*, or *Left atrium cavity* of the heart
   - [*class*] in the cardiac ultrasound
2. Include the chamber information
   - Left ventricular cavity in *two-chamber view* of the heart.
   - Left ventricular cavity in *two-chamber view* in the cardiac ultrasound.
3. Include the cycle
   - Left ventricular cavity in two-chamber view of the heart at the *end of the diastole cycle*.

- Left ventricular cavity in two-chamber view in the cardiac ultrasound at the *end of the diastole cycle*.

4. Include the gender

   - Left ventricular cavity in two-chamber view of the heart at the end of the diastole cycle of *a female*.
   - Left ventricular cavity in two-chamber view in the cardiac ultrasound at the end of the diastole cycle of *a female*.

5. Include the age

   - Left ventricular cavity in two-chamber view of the heart at the end of the diastole cycle of a *forty-six-year-old* female.
   - Left ventricular cavity in two-chamber view in the cardiac ultrasound at the end of the diastole cycle of a *forty-six-year-old* female.

6. Include the image quality

   - Left ventricular cavity in two-chamber view of the heart at the end of the diastole cycle of a 40-year-old female with *poor image quality*.
   - Left ventricular cavity in two-chamber view in the cardiac ultrasound at the end of the diastole cycle of a 40-year-old female with *poor image quality*.

7. Include the mask shape

   - Left ventricular cavity of *triangular shape* in two-chamber view of the heart at the end of the diastole cycle of a 40-year-old female with *poor image quality*.
   - Left ventricular cavity of *triangular shape* in two-chamber view in the cardiac ultrasound at the end of the diastole cycle of a 40-year-old female with *poor image quality*.

### G.2.3 Breast Ultrasound Images Dataset

The prompts for the Breast Ultrasound Images (BUSI) [1] dataset are listed below.

1. Presence of tumor

   - *[No] tumor* in the breast ultrasound image

2. Tumor Type

   - *Benign* tumor in the breast ultrasound image
   - *Regular-shaped* tumor in the breast ultrasound image

3. Tumor Number

   - *Two* benign tumors in the breast ultrasound image
   - *Two* regular-shaped tumors in the breast ultrasound image

4. Tumor Coverage

   - Two *medium* benign tumors in the breast ultrasound image
   - Two *medium* regular-shaped tumors in the breast ultrasound image

5. Tumor Location

   - Two medium benign tumors *at the center, left* in the breast ultrasound image
   - Two medium regular-shaped tumors *at the center, left* in the breast ultrasound image

6. Tumor Shape

   - Two medium *square-shaped* benign tumors at the center, left in the breast ultrasound image
   - Two medium *square-shaped* regular tumors at the center, left in the breast ultrasound image