

Computer Vision Mini Project Report

Abstract

In this report, we address a semantic segmentation task involving three classes: cats, dogs, and background. We develop and compare models based on encoder-decoder architectures, including U-Net, autoencoder-based variants, and CLIP-feature-enhanced designs. The best-performing model (U-Net) achieves an Intersection over Union (IoU) of 0.7992. Its robustness is evaluated under multiple perturbations. Additionally, the U-Net structure is adapted to support a prompt-based segmentation interface using points or boxes.

1. Introduction

This report addresses the task of semantic segmentation on the Oxford-IIIT Pet Dataset, focusing on classifying pixels into cats, dogs, and background. We implement and compare three segmentation models, U-Net, autoencoder-based, and CLIP-enhanced architectures, to identify the most effective approach in terms of accuracy, robustness, and adaptability.

To tackle dataset challenges such as class imbalance and limited diversity, we apply targeted preprocessing and augmentation. We also extend the system with a prompt-based segmentation interface, enabling user-guided refinement through points or boxes. The report presents the methodology, experiments, results, and an interactive user interface, providing a comprehensive evaluation of the strengths and limitations of each approach.

2. Data Preprocessing and Augmentations

2.1. Data Introduction

Overall review. This project utilizes a filtered subset of the Oxford-IIIT Pet Dataset. The dataset consists of two primary subsets: Train&Val (3,673 images) and Test (3,694 images). Each image is accompanied by a corresponding RGB segmentation mask. The segmentation task defines four pixel classes, each with a specific RGB colour mapping: **black** (0, 0, 0) represents the background, **white** (255, 255, 255) denotes uncertain re-

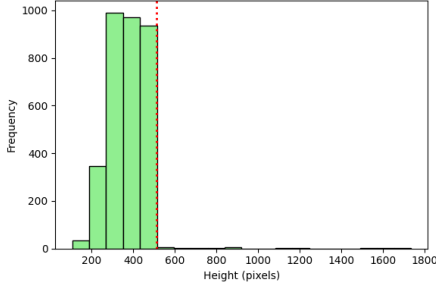
gions, **red** (128, 0, 0) corresponds to cats, and **green** (0, 128, 0) indicates dogs.

Uncertain regions, frequently located around cat boundaries and including elements such as clothing, accessories, or watermark overlays, are regarded as noise and hence treated as background during preprocessing. The RGB segmentation masks, initially unsuitable for loss calculations due to their 0–255 value range, are converted into discrete class maps, encoding the background (black and white regions) as class 0, cats as class 1, and dogs as class 2 (detailed in *preprocessing.py*, Appendix C.2). To prevent information leakage, the Test set remains entirely isolated from training and validation processes. The Train&Val subset is further divided into training and validation subsets in a 9:1 ratio (Appendix C.4), with validation solely used for model selection.

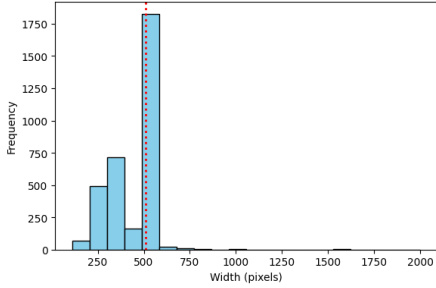
Reshape. Due to varying dimensions across images, batch operations and training efficiency are compromised. Consequently, the first preprocessing step standardizes all images and masks to a uniform resolution of (512×512) . Figure 1 demonstrates that only 1.27% images exceed this dimension, either width or height bigger than 512. It implies minimal information loss from resizing. The (512×512) resolution was specifically chosen to:

1. Ensure compatibility with standard convolutional neural network architectures (e.g., U-Net), which commonly use square-shaped inputs.
2. Facilitate valid feature map sizes through the downsampling and upsampling stages typical in encoder-decoder segmentation models, by adopting dimensions divisible by 16.

Resizing is performed using bilinear interpolation for images to maintain visual consistency, while nearest-neighbor interpolation is used for masks to preserve discrete class labels accurately. Despite these preparations, the dataset remains limited in both size and class balance — with a significant skew toward dog images and fewer variations in appearance or lighting. These issues motivate the need for targeted data augmentation strategies, detailed in the following subsection.



(a) Histogram of image heights in the raw training set. The red dotted line indicates the 512-pixel re-shape target.



(b) Histogram of image widths in the raw training set. The red dotted line indicates the 512-pixel re-shape target.

Figure 1. Distributions of image dimensions in the raw training dataset. A threshold of 512 pixels is highlighted to evaluate how many images exceed standard input size constraints.

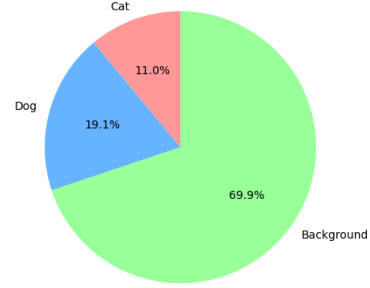
2.2. Data Augmentation

To address the dataset limitations discussed in Section 2.1, this augmentation pipeline serves two primary goals: 1) mitigating class imbalance; and 2) expanding the training dataset to improve model accuracy, generalisation, and robustness.

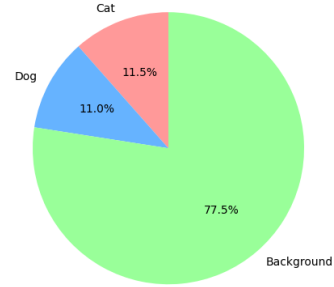
Class Imbalance Handling. The original dataset exhibited significant imbalance (Figure 2a), with dogs dominating both image and pixel counts. To counter this, we applied heavier augmentation to cat images—doubling their count relative to dogs—resulting in a more balanced dataset (Figure 2b).

Augmentation Strategy. To simulate real-world variability, we applied spatial and photometric augmentations that preserved mask alignment. These include: 1) horizontal flips, 2) rotations, translations, and scaling, 3) brightness and contrast adjustments, 4) color jittering, and 5) sparse Gaussian blur. They simulate (examples seen in Figure 3) real-world variations and common image defects that a model may encounter in deployment to ensure the generalization and robustness of trained model.

These transformations were applied with moderate

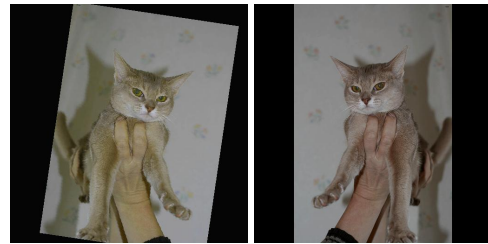


(a) Pixel class distribution of raw training dataset. Image-wise, there are 1055 cat images and 2250 dog images



(b) Pixel class distribution of augmented dataset. Image-wise, there are 8440 cat images and 9000 dog images

Figure 2. Pixel-wise class distributions of raw and augmented training datasets.



(a) Jittering

(b) Brightness adjustment



(c) Rotation + translation

(d) Horizontal flip

Figure 3. Examples of data augmentations applied to the same image. Note that multiple augmentation techniques are often applied together; the subcaptions indicate the most representative transformation in each case.

probability, ensuring that each augmented sample retained core semantic structure while introducing meaningful variation. All operations were carefully configured to avoid over-distortion or label corruption, particularly important in pixel-wise segmentation tasks where precise spatial correspondence between image and mask must be preserved.

Although the image labels are converted to class maps with values of 0, 1, and 2 (Section 2.1), the input colour images initially have pixel values in the range of 0 to 255. This wide range is not ideal for training deep neural networks, as it may lead to unstable gradients and slower convergence. Therefore, all input images are normalized to the range [0,1], improving numerical stability and facilitating more effective learning.

Outcome. Through this augmentation pipeline, the number of training images increased from 3,305 to 17,440. The class distribution became more balanced, and the diversity of training examples was significantly enhanced. This not only reduces the risk of overfitting but also equips the model to perform better in real-world conditions, especially on underrepresented cat images and in visually noisy or inconsistent scenes.

3. Segmentation networks design and implementation

3.1. Model Size and Selection

For an effective and fair comparison of segmentation performance across various architectures (see Section 4), all models were constrained to a similar parameter count of approximately 31×10^6 . A model size around 31 million parameters is commonly recognized in the literature as a balanced choice, effectively striking a trade-off between computational efficiency, memory usage, and segmentation accuracy, making it suitable for practical deployment on typical GPU hardware (He et al., 2016; Ronneberger et al., 2015). Specifically, this scale provides sufficient representational capacity without excessive computational demands, thus facilitating thorough experimentation within resource-limited settings typical for academic and applied research.

To achieve consistency in model comparisons, the parameter sizes of the U-Net, autoencoder-based, and CLIP-based segmentation models were explicitly matched. This uniformity ensures that observed performance differences primarily reflect architectural effectiveness rather than capacity disparities. The U-Net architecture, which requires training all parameters, establishes a baseline parameter count. In contrast, autoencoder and CLIP-based models incorporate pre-trained, frozen encoders, thereby significantly reducing the number of parameters that require training while maintaining comparable total model sizes.

We specifically selected ResNet-50 (RN50) as the visual backbone for the CLIP-based model due to several critical advantages. Compared to Vision Transformer (ViT) architectures like ViT-B/32, ResNet-50 employs smaller convolutional kernels (3×3), offering enhanced spatial precision essential for pixel-wise segmentation tasks. Moreover, RN50’s visual encoder consists of approximately 23.53×10^6 parameters, enabling adequate capacity for decoder components within the fixed parameter budget, thus aligning with the controlled experiment criteria established. Detailed parameter distributions for each model variant are summarized in Table 1.

Model	Config Size (MB)	Params #
U-Net	118.42	31,043,651
Autoencoder	119.59	31,350,275
CLIP	116.89	30,642,915

Table 1. Comparison of models based on configuration size and parameter count

3.2. Model Architecture

3.2.1 U-Net-based End-to-End Neural Network

The implemented U-Net segmentation model follows the established encoder-decoder architecture proposed by Ronneberger et al. [6]. Specifically adapted for RGB image segmentation, the network accommodates three input channels corresponding to preprocessed colour images (see Section 2.2), predicting pixel-level class maps during training. The encoder compresses input images into latent features with 1024 channels, using convolutional blocks followed by max-pooling operations (Figure 4). Each convolutional block contains two convolutional layers, batch normalization (BN), and rectified linear unit (ReLU) activation functions. BN stabilizes the learning process by normalizing inputs to each layer, which reduces internal covariate shift and accelerates convergence [2]. The ReLU activation function is used due to its computational efficiency, reduced risk of vanishing gradients, and effectiveness in promoting sparsity in activations, which aids generalization [3]. The decoder upsamples the latent features symmetrically using transpose convolutional layers. Skip connections between corresponding encoder and decoder blocks (grey arrows in Figure 4) mitigate gradient vanishing and enable the model to leverage both deep and shallow features effectively [6]. The final segmentation output is produced through a convolutional layer mapping to the segmentation class channels (Appendix D.1).

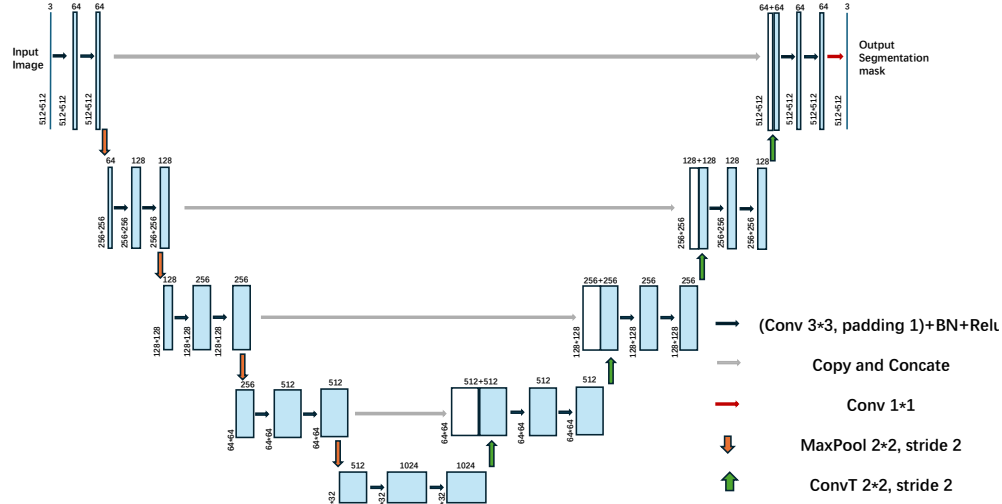


Figure 4. U-Net architecture used for baseline segmentation.

3.2.2 Autoencoder Pre-training Segmentation

The autoencoder-based pre-training model capitalizes on unsupervised feature extraction by initially training the network to reconstruct input images, approximating the identity function [4]. This unsupervised learning phase employs the mean squared error (MSE) loss, effectively encouraging the encoder to capture robust latent representations of the original data. The encoder sequentially reduces spatial dimensions using convolutional layers with strides, combined with BN and ReLU activations for stability and non-linearity. The decoder symmetrically reconstructs images using transpose convolutional layers (top row, Figure 5). Post pre-training, the encoder’s weights are frozen, serving as a feature extractor in a supervised segmentation network. A dedicated segmentation decoder, consisting of transpose convolutional layers, directly translates pre-trained latent representations into segmentation masks (Appendix D.2). This approach aligns with previous research demonstrating that unsupervised pre-training significantly improves downstream supervised segmentation performance by providing effective initialization and richer feature representations [1, 7].

3.2.3 CLIP-based Segmentation

The CLIP-based segmentation model integrates pre-trained visual representations from CLIP (Contrastive Language-Image Pre-training) [5], exploiting task-agnostic, semantically rich features extracted from its visual backbone. The weights of CLIP’s visual encoder are maintained frozen to preserve general feature representations beneficial for downstream tasks. These fea-

tures are then processed through a lightweight decoder designed to progressively recover the original image resolution using transpose convolutions and bilinear interpolation, thereby enhancing the spatial precision of the segmentation maps (Figure 6).

This architectural choice is supported by prior research indicating that leveraging pre-trained, general-purpose representations (such as those from CLIP) significantly improves semantic segmentation tasks by incorporating broad contextual and semantic insights [5, 9]. The use of transpose convolutions and bilinear interpolation strategically balances computational efficiency and output quality, ensuring accurate segmentation while maintaining manageable computational demands (Appendix D.3).

4. Experiments, Evaluation and Comparison

4.1. Experimental Settings

Data. We use the filtered dataset described in Sections 2.1 and 2.2, which contains RGB images of cats and dogs with pixel-wise segmentation masks. The original TrainVal set is split into training and validation subsets in a 9:1 ratio, detailed in Appendix C.4. All images and masks are resized to 512×512 , and the input images are normalized to the range $[0, 1]$. During training, extensive augmentations, including flipping, rotation, scaling, brightness and contrast adjustments, and Gaussian blurring, are applied. To mitigate class imbalance, the number of augmented cat images is doubled, resulting in a final training set of 17,440 images.

Models. The candidate segmentation models include

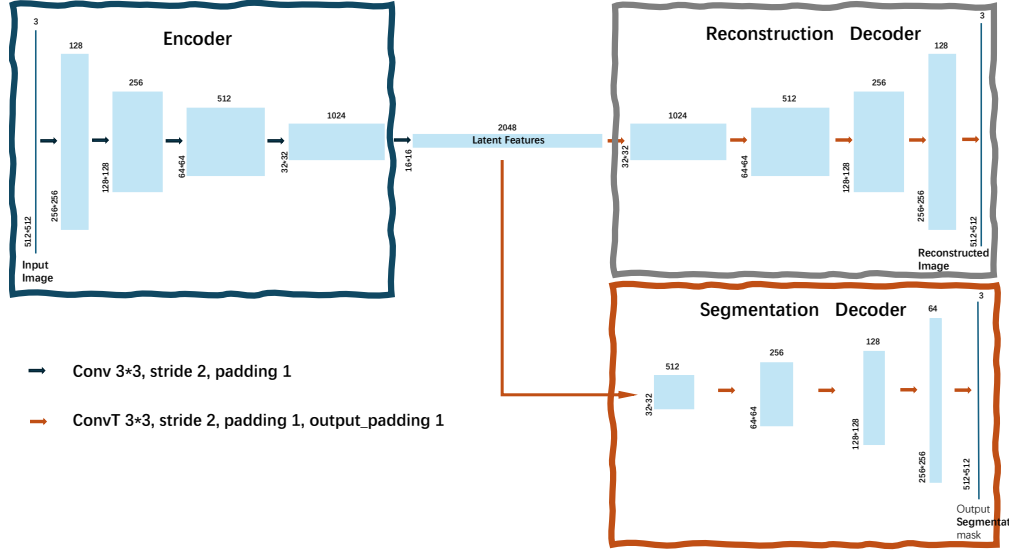


Figure 5. Autoencoder-based segmentation model with a frozen encoder.

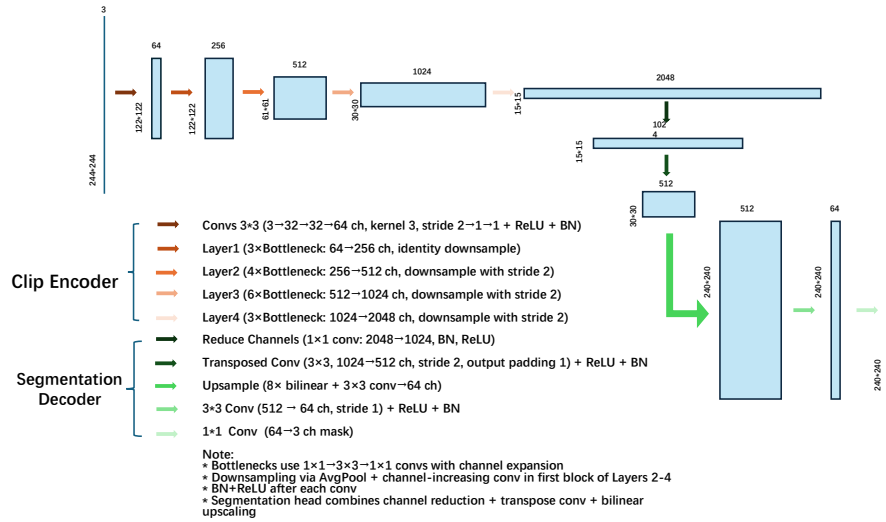


Figure 6. CLIP-enhanced model incorporating prompt-based feature guidance.

U-Net, Autoencoder, and CLIP-based architectures, as described in Section 3. For a fair comparison, the total number of parameters—including both trainable and frozen components—is controlled to be approximately 31 million for all models, as detailed in Section 3.1.

Hyperparameters. For fair comparison, we fix the learning rate to 1×10^{-4} , batch size to 4, epoch to 30 and use the Adam optimizer with a fixed random seed of 42. Each model checkpoint’s selection is based on cross-validation using the mean IoU, where historical checkpoints are evaluated to select the one with the best validation performance. All training and validation pro-

cesses are done on one GPU-4070s.

Loss functions. To effectively supervise this 3-class segmentation task, we evaluated multiple loss functions, including cross-entropy, Dice Loss, and IoU Loss. Ultimately, a composite of Focal Loss and Dice Loss was selected to jointly address class imbalance and region-level segmentation accuracy. Given predicted logits $\mathbf{z}_i \in \mathbb{R}^C$ at each pixel i , the *Focal Loss* is defined below and coded in Appendix E.3

$$\mathcal{L}_{\text{Focal}} = -\frac{1}{N} \sum_{i=1}^N \alpha_{y_i} \cdot (1 - \hat{p}_{i,y_i})^\gamma \cdot \log(\hat{p}_{i,y_i}) \quad (1)$$

where:

- N is the total number of pixels,
- $y_i \in \{0, 1, 2\}$ is the i th ground truth class index,
- $\hat{p}_{i,c} = \frac{\exp(z_{i,c})}{\sum_{k=1}^C \exp(z_{i,k})}$ is the softmax probability for class c at pixel i ,
- \hat{p}_{i,y_i} is the probability assigned to the true class,
- γ is the focusing parameter,
- α_{y_i} is an optional weighting factor for class y_i (can be scalar or class-dependent).

This formulation focuses on training on misclassified pixels by down-weighting the loss contribution of well-classified examples. The loss is averaged over all pixels.

The *Dice Loss* is defined as:

$$\mathcal{L}_{\text{Dice}} = 1 - \frac{2 \sum_i p_i y_i + \epsilon}{\sum_i p_i + \sum_i y_i + \epsilon} \quad (2)$$

where p_i and y_i are the predicted and ground truth values respectively, and ϵ is a small constant to avoid division by zero.

The **combined loss function** is:

$$\mathcal{L} = (1 - \alpha) \cdot \mathcal{L}_{\text{Focal}} + \alpha \cdot \mathcal{L}_{\text{Dice}} \quad (3)$$

where $\alpha \in [0, 1]$ is a tunable parameter that controls the balance between the two components.

Note: Intersection-over-Union (IoU) was used as the validation metric to align with the baseline evaluation protocol.

Training. For the U-Net model, we apply full end-to-end training, where all parameters are learned from scratch with combined loss function as defined in Equation 3.

For the autoencoder-based model, training is divided into two phases. During the pre-training phase, the encoder is trained as a reconstruction network using the mean squared error (MSE) loss:

$$\mathcal{L}_{\text{MSE}} = \frac{1}{N} \sum_{i=1}^N \|x_i - \hat{x}_i\|^2 \quad (4)$$

where x_i and \hat{x}_i denote the original and reconstructed images, respectively. After pre-training, the encoder weights are frozen, and a segmentation decoder is trained on top using the same segmentation loss as the U-Net model (Equation 3).

In the CLIP-based model, the visual encoder (CLIP-RN50) remains frozen throughout training, while a lightweight decoder is trained to map high-level CLIP features to segmentation masks with the same segmentation loss function in Equation 3.

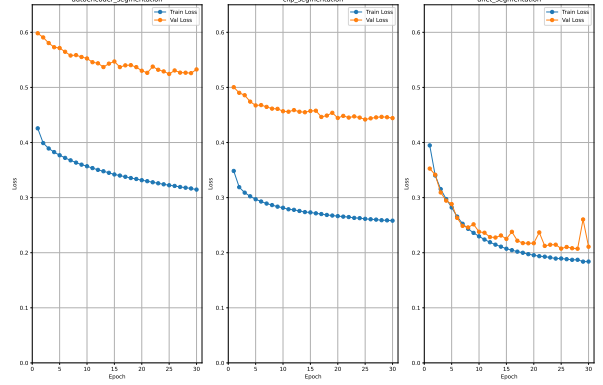


Figure 7. Curves for training loss(3) and validation loss (2). Plots generated by codes in Appendix E.2.

Validation and Model Selection. During training, we apply cross-validation with a fixed train-validation split ratio of 9:1. To select the best-performing checkpoints, we use the validation loss defined as the pure IoU loss, computed as $(1 - \text{IoU})$:

$$\mathcal{L}_{\text{val}} = 1 - \text{IoU} \quad (5)$$

This choice aligns with the final evaluation metric, ensuring consistency when comparing the three models.

After training, the model achieving the highest validation IoU score is selected as the final model for each architecture. The IoU is calculated as:

$$\text{IoU} = \frac{TP}{TP + FP + FN} \quad (6)$$

where TP , FP , and FN represent the number of true positive, false positive, and false negative pixels, respectively.

4.2. Quantitative Results and Analysis

Table 2 summarizes the segmentation performance of all models on the held-out validation set. The mean Intersection-over-Union (IoU) is reported, along with per-class IoU scores for background, cats, and dogs.

Model	Background IoU	Cats IoU	Dogs IoU	Mean IoU
U-Net	0.9315	0.7234	0.7429	0.7992
CLIP-based	0.8775	0.4864	0.5394	0.6344
Autoencoder	0.8003	0.2849	0.2171	0.4341

Table 2. Segmentation performance comparison on the validation set (per-class IoU and mean IoU).

Training Dynamics. Figure 7 shows the training and validation loss curves for the UNet, CLIP-based, and autoencoder segmentation models over 30 epochs. Among the three, UNet demonstrates the most efficient and effective training dynamics. It begins with a higher initial training loss (0.40) compared to CLIP-based (0.35)

and autoencoder (0.42), but exhibits a steep and consistent decline in both training and validation losses early in training. This rapid convergence indicates that the UNet architecture is well-suited to the segmentation task, likely due to its strong spatial inductive bias. In contrast, the CLIP-based model improves more slowly and shows a larger gap between training and validation losses, suggesting a tendency to overfit and limited alignment with pixel-level segmentation tasks. The autoencoder performs the worst, with the highest initial loss and the smallest reduction over time, indicating weak generalization and poor suitability for this task.

Overall, the training curves further support the quantitative results in Table 2, showing that U-Net not only achieves the best final IoU but also exhibits the most stable and effective training behavior.

U-Net. The U-Net model achieves the highest overall performance with a mean IoU of 0.7992. It significantly outperforms the other models across all classes, particularly in segmenting cats (0.7234) and dogs (0.7429). This result demonstrates the benefit of end-to-end training and its encoder-decoder architecture with skip connections, effectively preserving spatial details crucial for pixel-level segmentation. As shown in Figure 4, U-Net produces segmentations with well-defined edges and clear separation between foreground and background. However, a notable failure case is its misclassification of the chair beneath the cat as part of the dog, which suggests a limitation in distinguishing object context and inter-class relationships, especially when objects are physically close or visually similar. This may point to U-Net’s reliance on local texture and shape cues rather than broader semantic understanding.

CLIP-Featured. The CLIP-based model achieves a mean IoU of 0.6344, inferior from the performance of the U-Net. The model performs reasonably well on cats and dogs (0.4864 and 0.5394, respectively), highlighting the value of using pre-trained visual features. However, performance still lags behind U-Net, likely due to CLIP’s relatively coarse patch-wise representation and limited spatial granularity. These characteristics make it more difficult to capture object boundaries precisely. As shown in Figure 6, the CLIP-based model identifies most of the cat’s region but struggles with accurately delineating edges. This imprecision at boundaries may lead to merging or loss of fine structural information, particularly around thin or small regions like tails or legs.

Autoencoder. The autoencoder-based model shows the weakest performance, with a mean IoU of only 0.4341, and particularly poor results on cats (0.2849) and dogs (0.2171). These numbers reflect the model’s limited ability to capture task-specific features necessary for semantic segmentation. The encoder, pre-trained for

image reconstruction, likely focuses on global appearance and texture rather than precise object delineation, which hinders its utility for dense labeling. As shown in Figure 5, reconstructions from the autoencoder appear blurred, lacking critical spatial details, and resulting in fragmented, blocky segmentation outputs. Such fragmentation can be directly attributed to the loss of spatial resolution in the latent representations, compounded by the absence of architectural features like skip connections, essential for preserving high-frequency spatial information. Consequently, segmentation masks produced by this model are inconsistent, patchy, and inadequate for accurately segmenting smaller or complex-shaped regions.

5. Robustness Exploration

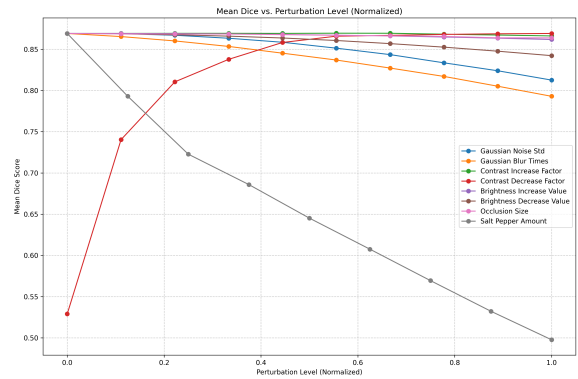


Figure 9. Mean Dice score across all perturbations

We further explore the robustness of the best-performing model (selected based on the validation dataset in Section 4) by applying eight types of perturbations to the test set:

1. **Gaussian pixel noise:** Zero-mean Gaussian noise with varying standard deviations is added independently to each pixel channel. The noise is sampled from a normal distribution $\mathcal{N}(0, \sigma^2)$ and clipped to remain within the valid image range $[0, 255]$.
2. **Gaussian blurring:** A fixed 3×3 Gaussian kernel is convolved multiple times with the image. The number of convolution iterations determines the blurring strength.
3. **Image contrast increase:** Pixel intensities are multiplied by a factor $\alpha > 1$, effectively stretching the dynamic range. The result is clipped to $[0, 255]$ to remain a valid image.
4. **Image contrast decrease:** Similar to the increase operation, but with a factor $0 < \alpha < 1$, reducing

the dynamic range and making the image appear more muted.

5. **Image brightness increase:** A fixed value is added to all pixel channels, uniformly brightening the image. Clipping ensures that pixel values remain within $[0, 255]$.
6. **Image brightness decrease:** A fixed value is subtracted from all pixel channels, uniformly darkening the image while also applying clipping.
7. **Image occlusion:** A square patch of specified size is randomly placed within the image and replaced with black pixels $(0, 0, 0)$, simulating partial obstruction.
8. **Salt-and-pepper noise:** A proportion of pixels is randomly replaced with either black $(0, 0, 0)$ or white $(255, 255, 255)$, simulating impulse noise.

These perturbations are applied at multiple severity levels to generate a range of distorted test datasets. The performance degradation curves under these perturbations are presented in Figures 9 and 10.

The model demonstrates strong robustness to brightness increase, contrast enhancement, and occlusion, showing only a slight Dice score drop (0.02 from 0.875) across perturbation levels.

In contrast, brightness reduction, Gaussian blur, and Gaussian noise cause exponential performance degradation, lowering the Dice score from 0.87 to 0.80. Contrast decrease and salt-and-pepper noise lead to the most severe drops: the former declines exponentially from 0.87 to 0.53, while the latter shows a linear decline to 0.50.

To explore worst-case scenarios, we assess one large object (dog) and one small object (cat) from the perturbed test set under the highest severity level (Figure 11).

For the dog, predictions remain close to the ground truth under benign perturbations (brightness/contrast increase, occlusion), indicating robustness—likely due to the large spatial context aiding segmentation even when corrupted.

In contrast, the cat example—representing smaller, low-contrast regions—suffers more. Brightness/contrast increase causes misclassification (e.g., the white cat’s black tail blends with the background), and occlusion hides critical regions. This suggests the model heavily relies on local texture and contrast.

With brightness reduction, noise, and blur, the dog’s predictions lose boundary precision, showing the importance of edge information. Interestingly, brightness reduction improves contrast for the cat, slightly aiding performance.

Gaussian blur severely impacts both examples by erasing fine details, leading to boundary erosion and shape distortion. Contrast decrease and salt-and-pepper noise result in sparse, fragmented predictions, highlighting the model’s reliance on intensity differences to delineate objects.

6. Prompt-based segmentation

Based on the best model architecture UNet as we discussed in the previous section, we develop a prompt-based segmentation model and build a UI for users to interact with their uploaded images.

6.1. Model Architecture

A prompt-based segmentation approach was implemented using a custom-designed point prompt encoder, which encodes spatial prompts (heatmaps) into feature embeddings. This encoder incorporates convolutional layers to transform prompt heatmaps into spatial representations, which are then concatenated with image features extracted from a simple convolutional backbone. The combined feature maps are processed through a decoder comprising convolutional layers to generate the segmentation masks. This approach allows for guided segmentation, leveraging external prompt information to improve segmentation accuracy.

6.2. Training

Data with Prompt. The data preprocessing involved generating spatial prompts (heatmaps) based on annotated points or bounding boxes provided in prompt files. During preprocessing, points or bounding boxes were randomly sampled from mask labels, balancing between foreground and background classes. These sampled points were then converted into Gaussian heatmaps with a specified sigma value. When bounding boxes were used, the heatmap consisted of a rectangular region corresponding to the box area. The final processed dataset thus consisted of RGB images, corresponding segmentation masks, and spatial prompt heatmaps to guide the segmentation model.

Loss Function. To train the model, we used the binary cross-entropy (BCE) loss, which is suitable for binary segmentation tasks. It compares the predicted probability \hat{y}_i for each pixel with the ground truth label y_i , penalizing incorrect predictions. The loss is defined as:

$$\mathcal{L}_{\text{BCE}} = -\frac{1}{N} \sum_{i=1}^N [y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)] \quad (7)$$

where N is the total number of pixels, $y_i \in \{0, 1\}$ is the ground truth label, and $\hat{y}_i \in [0, 1]$ is the predicted probability of the foreground class.

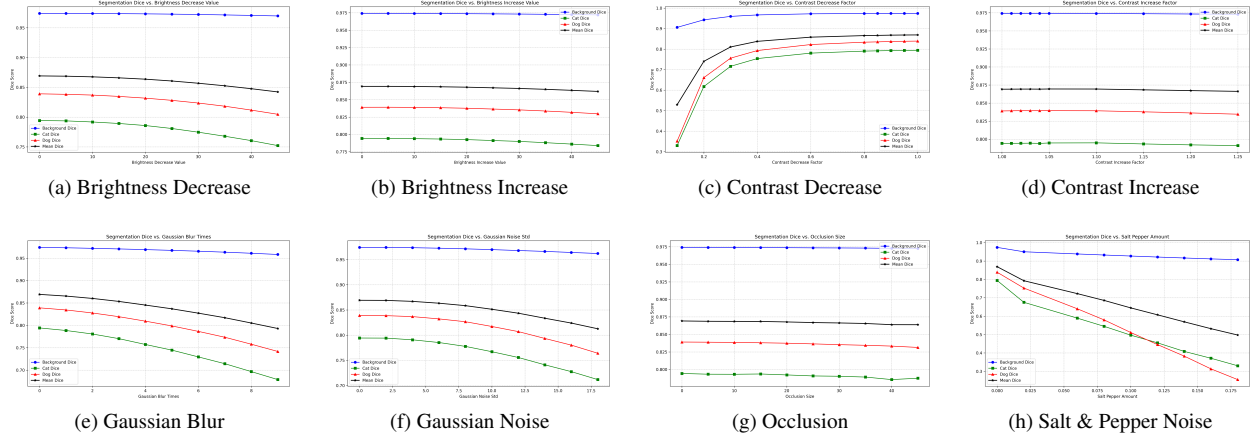


Figure 10. Dice scores under individual perturbation types

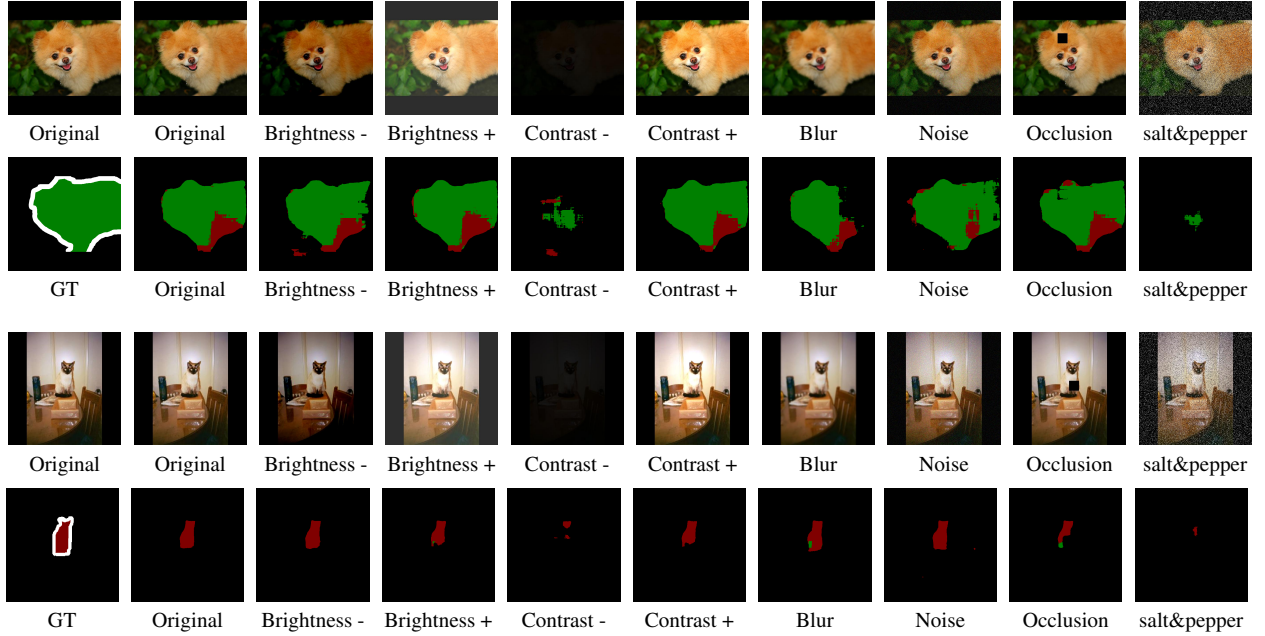


Figure 11. Visual comparison of model robustness under different perturbations. Each column corresponds to a perturbation type at its highest severity level: **Brightness**: pixel value shift (± 45), **Contrast**: multiplicative scaling ($\times 0.10 / \times 1.25$), **Gaussian Blur**: repeated convolution ($\times 9$), **Gaussian Noise**: standard deviation $\sigma = 18$, **Occlusion**: 45×45 black square, **Salt & Pepper Noise**: density = 0.18, Rows 1 & 2: dog input and prediction; Rows 3 & 4: cat input and prediction; Ground truths are shown in the first column of Rows 2 & 4.

Training Settings. The training was conducted using a custom PyTorch training loop specifically designed for prompt-based segmentation. The model input includes two components: RGB images, and generated prompt heatmaps. The output of the model is the predicted segmentation mask. Training was performed using early stopping based on validation loss, with parameters optimized via Adam or a similar optimizer, leveraging GPU acceleration for efficient computation.

6.3. Test

We achieve a mean IoU of 0.8780 on the test set using randomly sampled point and box prompts. The model demonstrates strong generalization on unseen samples. Figure 12 shows visual results from the test set, including input images, heatmaps of prompts, predicted masks, and ground truth masks. The model effectively segments foreground and background regions using either point or box prompts, showing especially strong performance on background separation. However, the

Table 3. Mean IoU / Dice under different perturbation types at various unified levels

Level	Gaussian Noise	Gaussian Blur	Contrast+	Contrast-	Brightness+	Brightness-	Occlusion	Salt & Pepper
1	0.7771/0.8692	0.7771/0.8692	0.7771/0.8692	0.7771/0.8692	0.7771/0.8692	0.7771/0.8692	0.7771/0.8692	0.7771/0.8692
2	0.7768/0.8690	0.7716/0.8654	0.7772/0.8693	0.7765/0.8688	0.7771/0.8692	0.7764/0.8687	0.7764/0.8687	0.6731/0.7930
3	0.7739/0.8670	0.7640/0.8601	0.7773/0.8694	0.7754/0.8680	0.7768/0.8690	0.7748/0.8676	0.7762/0.8686	0.6325/0.7596
4	0.7686/0.8634	0.7545/0.8535	0.7773/0.8694	0.7738/0.8669	0.7762/0.8686	0.7724/0.8660	0.7761/0.8685	0.5909/0.7227
5	0.7617/0.8585	0.7431/0.8453	0.7772/0.8693	0.7723/0.8659	0.7753/0.8680	0.7692/0.8637	0.7750/0.8678	0.5523/0.6859
6	0.7515/0.8514	0.7317/0.8370	0.7776/0.8695	0.7613/0.8583	0.7740/0.8671	0.7649/0.8607	0.7737/0.8669	0.5128/0.6452
7	0.7405/0.8435	0.7186/0.8272	0.7774/0.8694	0.7325/0.8378	0.7727/0.8662	0.7595/0.8569	0.7729/0.8664	0.4787/0.6075
8	0.7270/0.8336	0.7052/0.8171	0.7758/0.8683	0.6959/0.8105	0.7708/0.8649	0.7535/0.8526	0.7718/0.8656	0.4465/0.5694
9	0.7140/0.8240	0.6899/0.8052	0.7742/0.8673	0.6106/0.7404	0.7687/0.8635	0.7467/0.8477	0.7694/0.8639	0.4171/0.5322
10	0.6992/0.8126	0.6745/0.7930	0.7726/0.8662	0.4128/0.5289	0.7664/0.8619	0.7392/0.8423	0.7692/0.8638	0.3915/0.4976

model still exhibits instability under box prompts: edge predictions tend to be less accurate, and minor background regions are occasionally misclassified as foreground. This might be caused by the coarse nature of box prompts, which provide less precise spatial guidance compared to point annotations.

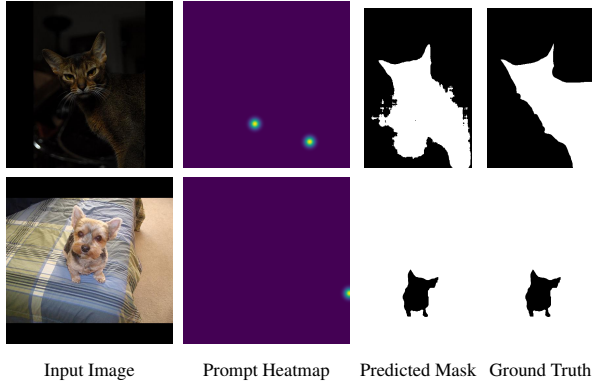


Figure 12. Qualitative results on the test set using point and box prompts.

6.4. UI Implementation

The UI is built with PyQt5, providing an interactive interface for prompt-based image segmentation. Users can load images, annotate with multiple points and boxes to guide segmentation using mouse, and adjust the threshold via a slider. The interface includes a central canvas for visualization and a control panel with buttons for loading images, running segmentation, clearing prompts, and saving masks. The system overlays segmentation results on the image and allows toggling the view. It supports combining multiple prompts for refined segmentation in a single session, offering a smooth and userfriendly experience. The GIFs in the supplemental materials show the end-to-end user workflow, including multi-points, boxes prompts, model inference, and visual feedback.

7. Limitation and Conclusion

In this project, we addressed the task of pixel-wise semantic segmentation for cats, dogs, and background using a filtered version of the Oxford-IIIT Pet Dataset. We developed and compared three architectures under a unified parameter budget: a classic U-Net, an autoencoder-based model with frozen encoders, and a CLIP-enhanced segmentation model leveraging pre-trained vision-language representations.

Among the three, the U-Net achieved the best performance, reaching a mean Intersection over Union (mIoU) of 0.7992. Its end-to-end training capability and skip connections allowed effective fusion of spatial details from early layers with deep semantic understanding, resulting in accurate segmentation even around complex object boundaries such as fur and ears.

In contrast, the CLIP-based model underperformed, highlighting several limitations. Although CLIP’s ResNet-50 encoder provides semantically rich features, it was originally trained for image-text alignment rather than dense prediction. Consequently, the extracted features lack the fine spatial granularity required for accurate segmentation. This challenge is exacerbated by the use of a *shallow decoder*, which lacks sufficient capacity to refine and upsample coarse features back to pixel-level precision. As a result, segmentation boundaries tend to be blurry or misaligned. Prior works such as MaskCLIP [10] and SegCLIP [8] also report similar issues and address them through deeper or multi-scale decoders.

The autoencoder-based model with frozen encoders faced similar issues: limited adaptability and a lack of spatial refinement reduced its effectiveness for detailed segmentation.

We applied targeted data augmentation to improve generalisation, address class imbalance, and expand the dataset to over 17,000 training samples. Robustness testing showed that U-Net maintained strong performance under most perturbations, although all models exhibited degradation under conditions like blur and low contrast.

Finally, we extended the U-Net with a prompt-based interface, supporting interactive segmentation via user-

provided points or bounding boxes. This extension demonstrates the model’s potential for real-world applications in interactive editing and annotation tools.

Future work may explore enhanced decoders for CLIP-based models, fine-tuning strategies for pre-trained encoders, multi-scale architectures like FPN or DeepLab, and domain adaptation for cross-species generalization.

References

- [1] Dumitru Erhan, Yoshua Bengio, Aaron Courville, Pierre-Antoine Manzagol, Pascal Vincent, and Samy Bengio. Why does unsupervised pre-training help deep learning? In *Journal of Machine Learning Research*, volume 11, pages 625–660, 2010. 4
- [2] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *Proceedings of the 32nd International Conference on Machine Learning (ICML)*, pages 448–456, 2015. 3
- [3] Vinod Nair and Geoffrey E. Hinton. Rectified linear units improve restricted boltzmann machines. In *Proceedings of the 27th International Conference on Machine Learning (ICML)*, pages 807–814, 2010. 3
- [4] Andrew Ng et al. Sparse autoencoder. *CS294A Lecture notes*, 72(2011):1–19, 2011. 4
- [5] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PmLR, 2021. 4
- [6] 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*, pages 234–241. Springer, 2015. 3
- [7] Pascal Vincent, Hugo Larochelle, Isabelle Lajoie, Yoshua Bengio, and Pierre-Antoine Manzagol. Stacked denoising autoencoders: Learning useful representations in a deep network with a local denoising criterion. In *Journal of Machine Learning Research*, volume 11, pages 3371–3408, 2010. 4
- [8] Yuchao Wang, Peihao Xu, Zhe Zhang, Xiang Xu, Yifan Wang, Yuhui He, and Xiaojuan Qi. Segclip: Leveraging clip for semantic segmentation. *arXiv preprint arXiv:2209.11888*, 2022. 10
- [9] Bowen Zhou, Junjie He, Xiang Yao, Philipp Krahenbuhl, Vladlen Koltun, and Alexei A. Efros. Image segmentation using text and image prompts. *arXiv preprint arXiv:2112.10003*, 2022. 4
- [10] Bowen Zhou, Chen Wang, Xinlong Wang, Dongdong Yu, Jifeng Dai, Yuning Lu, and Lu Yuan. Maskclip: Masked language-image pre-training for open-world semantic segmentation. *arXiv preprint arXiv:2208.08984*, 2022. 10

Appendix

A. Team Member Responsibilities

This appendix outlines the specific responsibilities and contributions of each team member. For tasks completed collaboratively, the individual portions are clearly specified.

Task	Team Member A	Team Member B
Data Preprocessing & Augmentation	Designed and implemented core preprocessing pipeline; generated balanced training datasets through augmentation	Developed preprocessing routines for prompt-based segmentation, including heatmap generation and spatial prompt encoding
Segmentation Model Development	Implemented and trained the autoencoder-based and CLIP-featured segmentation models	Implemented and trained the U-Net baseline and extended it for prompt-based segmentation
Evaluation & Analysis	Produced training plots and quantitative results for autoencoder and CLIP models	Conducted comparative evaluation and robustness analysis; visualized failure cases and strengths
Report Writing	Contributed to Sections 1–2 (Data preparation, Model design)	Contributed to Sections 3–4 (Evaluation, Robustness, Prompt-based extension)

Table 4. Breakdown of responsibilities by team member.

B. Code Structure Overview & Training Entry

B.1. Directory Tree

This appendix presents the structure and organization of the project source code. It outlines directories and scripts, highlighting their roles in dataset preprocessing, model architecture definitions, training logic, and inference routines. A central training script `main.py` provides a unified interface for training various segmentation models, including U-Net, AutoEncoder-based, CLIP-based, and prompt-enhanced variants. The modular design ensures clarity, flexibility, and ease of experimentation throughout the development pipeline.

```
src/
|-- data/
|   |-- __init__.py
|   |-- PetDataset.py           % Dataset definitions for all segmentation
|                               % and prompt models
|   |-- preprocess_data_with_prompt.py % Preprocess prompt-based datasets,
|                                   % taking CLI args
|   |-- preprocess_data.py         % Preprocess segmentation datasets,
|                                   % taking CLI args
|   |-- preprocessing.py           % General preprocessing functions
|                                   % (class mapping, reshape, augmentor)
|   |-- train_val_split.py         % Dataset splitting utility
|-- models/
|   |-- unet_segmentation.py       % U-Net model definition
|   |-- autoencoder_segmentation.py % AutoEncoder & Segmentation Head
|   |-- clip_segmentation.py       % CLIP-based segmentation model
|   |-- prompt_segmentation.py     % Prompt-enhanced U-Net
|-- utils/
|   |-- inference.py
|   |-- loss_functions.py          % Loss definitions (Focal, Dice, IoU)
```


	-- restore_image_size.py	% Restore predicted mask to original size
	-- training_plot.py	% Plot training/validation curves
	-- training.py	% Training loops (standard & prompt)
	-- compute_IoU.py	% Script to compute evaluation metric
	-- main.py	% Entry point for training with CLI args
	-- run_inference.py	% Perform model inference

B.2. Training Entry Script: main.py

The `main.py` script allows training four types of models via a unified CLI:

- **Mode 0:** U-Net-based segmentation
- **Mode 1:** AutoEncoder with optional pretraining
- **Mode 2:** CLIP-based segmentation (ResNet-50 backbone)
- **Mode 3:** Prompt-based segmentation using (x, y, class) inputs

The script handles argument parsing, dataset loading and splitting, model instantiation, training, and saving.

main.py

```

1 #!/usr/bin/env python
2
3 import torch
4 import argparse
5 import clip
6 import torch.optim as optim
7 from pathlib import Path
8 from sklearn.model_selection import train_test_split
9 from torch.utils.data import DataLoader
10 from data.preprocessing import *
11 from data.PetDataset import PetDataset, PetDatasetWithPrompt
12 # import unet related functions
13 from models.unet_segmentation import UNet
14 # import autoencoder related functions
15 from models.autoencoder_segmentation import (AutoEncoder, AutoEncoderSegmentation,
16                                              pretrain_autoencoder)
17 # import CLIP related functions
18 from models.clip_segmentation import CLIPSegmentationModel
19 # import prompt based model
20 from models.prompt_segmentation import PromptSegmentation
21 from utils.loss_functions import CombinedFocalDiceLoss, BinaryFocalLoss, IouLoss
22 from utils.training_plot import training_plot
23 from utils.training import training, prompt_training
24
25 # Define Arguments
26 parser = argparse.ArgumentParser(description="Train segmentation model")
27 parser.add_argument("--img_dir", type=str, default="Dataset/Processed/color", help="Directory_for_training_images")
28 parser.add_argument("--msk_dir", type=str, default="Dataset/Processed/label", help="Directory_for_training_masks")
29 parser.add_argument("--pnt_dir", type=str, default="Dataset/ProcessedWithPrompt/color/points", help="Directory_for_training_prompt_points")
30 parser.add_argument("--mode", type=int, choices=[0, 1, 2, 3], required=True, help="0_for_AutoEncoder, 1_for_CLIP-based_segmentation")
31 parser.add_argument("--pretrain", type=int, choices=[0, 1], default=0, help="1_to_pretrain_autoencoder")
32 parser.add_argument("--epochs", type=int, required=True, help="Number_of_training_epochs")
33 parser.add_argument("--batch_size", type=int, default=4, help="Batch_size_for_training")
34 parser.add_argument("--lr", type=float, default=1e-4, help="Learning_rate")
35 parser.add_argument("--save_dir", type=str, default="params/", help="Path_to_save_the_model")
36 parser.add_argument("--patience", type=int, default=5, help="Patience_for_early_stopping")
37 args = parser.parse_args()
38
39 # Define Device

```

```

40 device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
41 print(f'Training_run_on:{device}')
42 print(torch.cuda.get_device_name(0)) if torch.cuda.is_available() else print('no_cuda_devices')
43 # Ensure save directory exists
44 Path(args.save_dir).mkdir(parents=True, exist_ok=True)
45
46 # Load Dataset Paths
47 img_paths = sorted(Path(args.img_dir).glob("*.*)" )
48 msk_paths = sorted(Path(args.msk_dir).glob("*.*)" )
49 pnt_paths = sorted(Path(args.pnt_dir).glob("*.*)" )
50
51 if args.mode == 3:
52     train_img_paths, val_img_paths, \
53     train_msk_paths, val_msk_paths, \
54     train_pnt_paths, val_pnt_paths = train_test_split(img_paths, msk_paths, pnt_paths, test_size
55                                                         =0.2, random_state=42)
56 else:
57     train_img_paths, val_img_paths, \
58     train_msk_paths, val_msk_paths = train_test_split(img_paths, msk_paths, test_size=0.2,
59                                                         random_state=42)
60
61 # Select Model
62 #####
63 ##### U-Net #####
64 #####
65 if args.mode == 0:
66     print("Training_Unet-Based_Segmentation_Model...")
67     # define the model
68     model = UNet()
69
70     # define training data split
71     train_dataset = PetDataset(
72         img_paths=train_img_paths,
73         msk_paths=train_msk_paths,
74         transform=standard_transform)
75     val_dataset = PetDataset(
76         img_paths=val_img_paths,
77         msk_paths=val_msk_paths,
78         transform=standard_transform)
79
80     train_loader = DataLoader(train_dataset, batch_size=args.batch_size, shuffle=True,
81                               num_workers=4)
82     val_loader = DataLoader(val_dataset, batch_size=args.batch_size, shuffle=False, num_workers
83                             =4)
84
85     # define save paths
86     train_plot = "unet_segmentation.png"
87     training_plot_save_name = Path(args.save_dir) / train_plot
88     save_name = "unet_segmentation.pth"
89
90     # Define Loss Functions and Optimizer
91     criterion= CombinedFocalDiceLoss(alpha=0.5, gamma=2.0)
92
93     optimizer = optim.Adam(model.parameters(), lr=args.lr)
94
95     # Train Model
96     history = training(
97         model=model,
98         train_loader=train_loader,
99         val_loader=val_loader,
100         train_criterion=criterion,
101         val_criterion=criterion,
102         optimizer=optimizer,
103         num_epochs=args.epochs,
104         device=device,
105         save_dir=args.save_dir,
106         save_name=save_name,

```

```

104         patience=args.patience
105     )
106
107     training_plot(history,save_path=training_plot_save_name)
108
109     #####
110     ##### Autoencoder #####
111     #####
112 elif args.mode == 1:
113     # define training data split
114     train_dataset = PetDataset(
115         img_paths=train_img_paths,
116         msk_paths=train_msk_paths,
117         transform=standard_transform)
118     val_dataset = PetDataset(
119         img_paths=val_img_paths,
120         msk_paths=val_msk_paths,
121         transform=standard_transform)
122
123     train_loader = DataLoader(train_dataset, batch_size=args.batch_size, shuffle=True,
124                               num_workers=4)
125     val_loader = DataLoader(val_dataset, batch_size=args.batch_size, shuffle=False, num_workers
126                             =4)
127     # define config file name
128     pretrain_model_save_name = "train_autoencoder_pretrain_unified_size.pth"
129     model_save_name = "train_autoencoder_segmentation_unified_size.pth"
130     pretrain_plot = "train_autoencoder_pretrain_unified_size.png"
131     seg_train_plot = "train_autoencoder_segmentation_unified_size.png"
132 if args.pretrain == 1: # the pretrain mode
133     print('Autoencoder_Pretrain...')
134     autoencoder = AutoEncoder()
135     training_plot_save_name = Path(args.save_dir) / pretrain_plot
136     history = pretrain_autoencoder(autoencoder, train_loader, val_loader, num_epochs=args.
137                                     epochs,
138                                     save_dir=args.save_dir, save_name=pretrain_model_save_name
139                                     ,
140                                     device=device, patience=args.patience)
141     training_plot(history,save_path=training_plot_save_name)
142     print("Autoencoder_pretrained_and_saved.")
143 else: # the segmentation head train mode (encode freezed)
144     print('Training_Autoencoder-Based_Segmentation_Model...')
145     autoencoder = AutoEncoder()
146     pretrain_config_path = Path(args.save_dir) / pretrain_model_save_name
147     pretrain_state_dict = torch.load(pretrain_config_path, map_location=device)
148     autoencoder.load_state_dict(pretrain_state_dict)
149     training_plot_save_name = Path(args.save_dir) / seg_train_plot
150     model = AutoEncoderSegmentation(autoencoder.encoder, num_classes=3)
151     train_criterion = CombinedFocalDiceLoss()
152     val_criterion = IouLoss()
153     # filter(function,iterable) and lambda arguments: expression freeze the pretraining
154     # and train the decoder only
155     optimizer = optim.Adam(filter(lambda p: p.requires_grad, model.parameters()), lr=args.lr)
156     history = training(
157         model, train_loader, val_loader,
158         train_criterion=train_criterion, val_criterion=val_criterion,
159         optimizer=optimizer,
160         num_epochs=args.epochs, device=device,
161         save_dir=args.save_dir,
162         save_name=model_save_name,
163         patience=args.patience
164     )
165     training_plot(history,save_path=training_plot_save_name)
166     print("Autoencoder_pretrained_and_saved.")
167
168     #####
169     ##### CLIP #####
170     #####
171 elif args.mode == 2:

```

```

168 print("Training_CLIP-Based_Segmentation_Model...")
169 clip_model, _ = clip.load("RN50", device=device)
170 model = CLIPSegmentationModel(clip_model, num_classes=3)
171 train_dataset = PetDataset(
172     img_paths=train_img_paths,
173     msk_paths=train_msk_paths,
174     resize_fn=resize_with_padding,
175     resize_target_size=224,
176     transform=clip_transform)
177 val_dataset = PetDataset(
178     img_paths=val_img_paths,
179     msk_paths=val_msk_paths,
180     resize_fn=resize_with_padding,
181     resize_target_size=224,
182     transform=clip_transform)
183
184 train_loader = DataLoader(train_dataset, batch_size=args.batch_size, shuffle=True,
185     num_workers=4)
186 val_loader = DataLoader(val_dataset, batch_size=args.batch_size, shuffle=False, num_workers
187     =4)
188
189 # define save paths
190 train_plot = "train_clip_segmentation_unified_size.png"
191 training_plot_save_name = Path(args.save_dir) / train_plot
192 save_name = "train_clip_segmentation_unified_size.pth"
193
194 # Define Loss Functions and Optimizer
195 train_criterion = CombinedFocalDiceLoss()
196 val_criterion = IouLoss()
197
198 optimizer = optim.Adam(model.parameters(), lr=args.lr)
199
200 # Train Model
201 history = training(
202     model=model,
203     train_loader=train_loader,
204     val_loader=val_loader,
205     train_criterion=train_criterion,
206     val_criterion=val_criterion,
207     optimizer=optimizer,
208     num_epochs=args.epochs,
209     device=device,
210     save_dir=args.save_dir,
211     save_name=save_name,
212     patience=args.patience
213 )
214
215 training_plot(history, save_path=training_plot_save_name)
216
217 elif args.mode == 3:
218     print('Training_Prompt-based_Segmentation_model...')
219     # define model
220     model = PromptSegmentation(unet_in_channels=3,
221         prompt_dim=1024,
222         unet_init_features=64)
223
224     train_dataset = PetDatasetWithPrompt(
225         img_paths=train_img_paths,
226         msk_paths=train_msk_paths,
227         pnt_paths= train_pnt_paths,
228         transform=standard_transform)
229     val_dataset = PetDatasetWithPrompt(
230         img_paths=val_img_paths,
231         msk_paths=val_msk_paths,
232         pnt_paths = val_pnt_paths,
233         transform=standard_transform)
234
235     train_loader = DataLoader(train_dataset, batch_size=args.batch_size, shuffle=True,
236         num_workers=4)

```

```

233     val_loader = DataLoader(val_dataset, batch_size=args.batch_size, shuffle=False, num_workers
                             =4)
234
235     # define save paths
236     train_plot = "prompt_unet_segmentation.png"
237     training_plot_save_name = Path(args.save_dir) / train_plot
238     save_name = "prompt_unet_segmentation.pth"
239
240     # Define Loss Functions and Optimizer
241     criterion= BinaryFocalLoss()
242
243     optimizer = optim.Adam(model.parameters(), lr=args.lr)
244
245     # Train Model
246     history = prompt_training(
247         model=model,
248         train_loader=train_loader,
249         val_loader=val_loader,
250         train_criterion=criterion,
251         val_criterion=criterion,
252         optimizer=optimizer,
253         num_epochs=args.epochs,
254         device=device,
255         save_dir=args.save_dir,
256         save_name=save_name,
257         patience=args.patience
258     )
259
260     training_plot(history, save_path=training_plot_save_name)
261
262 else:
263     raise ValueError("Invalid_mode._Use_0_for_U-Net,_1_for_autoencoder-based_segmentation,_2_for_
                       CLIP-based_segmentation_or_3_for_prompt_segmentation.")

```

C. Data Preprocessing

This appendix provides the implementation details of the data preprocessing pipeline used in the segmentation project. It includes dataset loaders, preprocessing utilities, and command-line preprocessing scripts for both standard and prompt-based segmentation tasks. These components handle colour-class conversions, input resizing, data augmentation, and train-validation splits, ensuring consistent and reproducible input preparation across all training modes.

C.1. Torch Dataset Loader *PetDataset.py*

This file defines the dataset loaders for both standard segmentation and prompt-based segmentation tasks.

PetDataset.py

```

1 import torch
2 import random
3 import numpy as np
4 from PIL import Image
5 from torch.utils.data import Dataset
6 from data.preprocessing import color2class
7
8
9 class PetDataset(Dataset):
10     """
11     A PyTorch Dataset for loading and preprocessing image-mask pairs for segmentation.
12
13     Args:
14         image_paths (list of str): List of file paths to the images.
15         mask_paths (list of str): List of file paths to the corresponding masks.
16         resize_fn (callable, optional): A function to resize the image & mask.
17         resize_target_size (int): the target size of input images to the model. (assuming it's a
18             square image)
19         augment_fn (callable, optional): A function to apply data augmentation to the image.

```



```

19         transform (callable, optional): Function to apply final normalization (e.g., CLIP
20             transform).
21
22     Returns:
23         Tuple[torch.Tensor, torch.Tensor, Tuple[int, int], str]:
24             - image: A tensor of shape (C, H, W) representing the normalized image.
25             - mask: A tensor of shape (H, W), with values {0,1} or {0,2} (background vs cat or
26                 background vs dog).
27             - initial_img_size: Tuple (H, W) representing the original image size.
28             - img_path: str name of the color image
29
30     """
31     def __init__(self, img_paths, msk_paths, resize_fn=None, resize_target_size=None, transform=
32         None):
33         self.img_paths = img_paths
34         self.msk_paths = msk_paths
35         self.resize_fn = resize_fn # Function to resize images and masks differently
36         self.resize_target_size = resize_target_size
37         self.transform = transform
38         # make sure # of images is equal to # of masks
39         assert len(self.img_paths) == len(self.msk_paths), "Mismatch_between_images_and_masks."
40
41     def __len__(self):
42         return len(self.img_paths)
43
44     def __getitem__(self, idx):
45         # Load and convert image to RGB
46         img_path = str(self.img_paths[idx])
47         img = Image.open(img_path).convert("RGB")
48         img = np.array(img)
49         initial_img_size = tuple(img.shape[:2])
50
51         # Load and process mask
52         msk_path = str(self.msk_paths[idx])
53         msk = Image.open(msk_path).convert("RGB")
54         msk = np.array(msk) # Convert to NumPy array
55
56         # Convert color mask to class labels
57         msk = color2class(msk) # Convert mask from RGB to class labels
58
59         # Resize if a resizing function is provided
60         if self.resize_fn:
61             img = self.resize_fn(img, target_size=self.resize_target_size, is_mask=False)
62             msk = self.resize_fn(msk, target_size=self.resize_target_size, is_mask=True)
63
64         # Apply transform only to the image (not the mask)
65         if self.transform:
66             img = self.transform(Image.fromarray(img)) # only apply the normalization to img
67
68         # Convert mask to tensor
69         msk = torch.tensor(msk, dtype=torch.long)
70
71         return img, msk, initial_img_size, img_path
72
73 class PetDatasetWithPrompt(Dataset):
74     def __init__(self, img_paths, msk_paths, pnt_paths, resize_fn=None,
75         resize_target_size=None, transform=None, load_multiple_points=False,
76         use_box_prompt=True):
77         self.img_paths = img_paths
78         self.msk_paths = msk_paths
79         self.pnt_paths = pnt_paths
80         self.resize_fn = resize_fn
81         self.resize_target_size = resize_target_size
82         self.transform = transform
83         self.load_multiple_points = load_multiple_points
84         self.use_box_prompt = use_box_prompt
85
86         assert len(self.img_paths) == len(self.msk_paths) == len(self.pnt_paths), \
87             "Mismatch_between_images,_masks,_and_point_files."

```

```

83
84 def __len__(self):
85     return len(self.img_paths)
86
87 def load_points_from_file(self, point_path):
88     with open(point_path, 'r') as f:
89         line = f.readline().strip()
90
91     if not line:
92         raise ValueError(f"Empty_prompt_file:_{point_path}")
93
94     try:
95         parts = list(map(int, line.strip().split(',')))
96     except Exception as e:
97         raise ValueError(f"Failed_to_parse_prompt_file:_{point_path}_with_line:_{line}'_{e}({e})")
98
99     if len(parts) == 5:
100         # Box format: x1, y1, x2, y2, cls
101         x1, y1, x2, y2, cls = parts
102         return [], True, (x1, y1, x2, y2), cls
103     elif len(parts) == 3:
104         # Point format: x, y, cls
105         x, y, cls = parts
106         return [(x, y, cls)], False, None, -1
107     else:
108         raise ValueError(f"Invalid_prompt_format_in:_{point_path}_got_{len(parts)}_values")
109
110
111
112 def generate_gaussian_heatmap(self, H, W, x, y, sigma=10):
113     xs = np.arange(W)
114     ys = np.arange(H)
115     xs, ys = np.meshgrid(xs, ys)
116     g = np.exp(-((xs - x) ** 2 + (ys - y) ** 2) / (2 * sigma ** 2))
117     return g.astype(np.float32)
118
119 def __getitem__(self, idx):
120     img_path = str(self.img_paths[idx])
121     img = Image.open(img_path).convert("RGB")
122     img = np.array(img)
123     img_h, img_w = img.shape[:2]
124     initial_img_size = (img_h, img_w)
125
126     msk_path = str(self.msk_paths[idx])
127     msk = Image.open(msk_path).convert("RGB")
128     msk = np.array(msk)
129     msk = color2class(msk)
130
131
132     if self.resize_fn:
133         img = self.resize_fn(img, target_size=self.resize_target_size, is_mask=False)
134         msk = self.resize_fn(msk, target_size=self.resize_target_size, is_mask=True)
135
136     point_path = str(self.pnt_paths[idx])
137     points, is_box, box_coords, box_class = self.load_points_from_file(point_path)
138
139
140     H, W = msk.shape[:2]
141
142     if is_box and self.use_box_prompt:
143         x1, y1, x2, y2 = box_coords
144         heatmap = np.zeros((H, W), dtype=np.float32)
145         heatmap[y1:y2+1, x1:x2+1] = 1.0
146         point_class = box_class
147         # Use center of box as representative point
148         x = (x1 + x2) // 2
149         y = (y1 + y2) // 2

```

```

150     else:
151         if points == []:
152             print(is_box, box_coords, box_class)
153             raise ValueError(f"Unknown_prompt_format_in:_{point_path}_with_points_{}")
154
155         x, y, cls = points[0]
156         heatmap = self.generate_gaussian_heatmap(H, W, x, y, sigma=10)
157         point_class = cls
158
159     if self.transform:
160         img = self.transform(Image.fromarray(img))
161
162     output = {
163         'image': img,
164         'gt_mask': torch.tensor(msk, dtype=torch.long).unsqueeze(0),
165         'prompt_heatmap': torch.tensor(heatmap, dtype=torch.float32).unsqueeze(0),
166         'prompt_point': torch.tensor([x, y], dtype=torch.long),
167         'initial_img_size': initial_img_size,
168         'img_path': img_path
169     }
170
171     if point_class != -1:
172         output['point_class'] = torch.tensor(point_class, dtype=torch.long)
173         if point_class not in [0,1,2]:
174             print(f"invalid_point_class_in:_{point_path}")
175             if point_class==255:
176                 output['point_class'] = torch.tensor(0, dtype=torch.long)
177                 print("change_to_0")
178
179     return output

```

C.2. Preprocessing Utilities *preprocessing.py*

This file contains core utility functions for data preprocessing, including:

1. Mapping functions between color-encoded masks and class arrays.
2. A padding-based image resizing function for uniform input dimensions.
3. Image normalization transforms.
4. Data augmentation using the albumentations library.
5. A function to generate point-wise masks for the prompt-based model.

preprocessing.py

```

1  # All preprocessing utilities (image/mask transformations & augmentations)
2
3  import numpy as np
4  import cv2 as cv
5  import torchvision.transforms as transforms
6  from pathlib import Path
7  from PIL import Image
8  import torch
9  from tqdm import tqdm
10 import albumentations as A
11
12 # -----
13 # Color Mapping Utilities
14 # -----
15
16 def color2class(img: np.ndarray) -> np.ndarray:
17     """
18     Convert a label RGB image to pixel-wise class labels.
19
20     Parameters:

```

```

21     img (np.ndarray): Input image as a NumPy array (expected shape: (H, W, 3)).
22
23     Returns:
24         np.ndarray: Output array pixel-wise class labels (expected shape: (H, W)).
25     """
26
27     color_to_class = {
28         (0, 0, 0): 0,          # Black -> Class 0
29         (255, 255, 255): 0,    # White -> Class 0
30         (128, 0, 0): 1,        # Dark Red -> Class 1
31         (0, 128, 0): 2         # Green -> Class 2
32     }
33     h,w,_ = img.shape
34     class_map = np.zeros((h,w), dtype=np.uint8)
35     img_resaped = img.reshape(-1,3)
36     class_map_resaped = np.zeros(img_resaped.shape[0], dtype=np.uint8)
37
38     # Assign class labels
39     for color, class_id in color_to_class.items():
40         mask = np.all(img_resaped == color, axis=1)
41         class_map_resaped[mask] = class_id
42
43     # Reshape back to original dimensions
44     class_map = class_map_resaped.reshape(h, w)
45
46     return class_map
47
48 def class2color(class_map: np.ndarray) -> np.ndarray:
49     """
50     Convert a pixel-wise class label map to an RGB image representation.
51
52     Parameters:
53         class_map (np.ndarray): 2D numpy array of shape (H, W) containing class labels.
54
55     Returns:
56         np.ndarray: Output RGB image as a NumPy array of shape (H, W, 3).
57     """
58
59     # Define the mapping from class labels to RGB colors
60     class_to_color = {
61         0: (0, 0, 0),          # Class 0 -> Black (background)
62         1: (128, 0, 0),        # Class 1 -> Dark Red (cat)
63         2: (0, 128, 0)         # Class 2 -> Green (dog)
64     }
65
66     # Get image dimensions
67     h, w = class_map.shape
68
69     # Initialize an empty RGB image
70     color_image = np.zeros((h, w, 3), dtype=np.uint8)
71
72     # Assign colors based on class labels
73     for class_id, color in class_to_color.items():
74         mask = (class_map == class_id)
75         color_image[mask] = color # Assign the corresponding RGB color
76
77     return color_image
78
79 # -----
80 # Resizing Utilities
81 # -----
82
83 def resize_with_padding(img: np.ndarray, target_size, fill=0, is_mask=False) -> np.ndarray:
84     """
85     Resize an image while maintaining its aspect ratio and pad it to a square.
86
87     Args:
88         img (np.ndarray): Input image (H, W, C) or (H, W) if grayscale/mask.

```

```

89     target_size (int, optional): The target width and height (default: 224).
90     fill (int or tuple, optional): Padding color, either an int (grayscale) or
91         (R, G, B) tuple for color images. Default is black (0).
92     is_mask (bool, optional): If True, uses NEAREST interpolation for masks
93         to avoid artifacts. Default is False (for normal images).
94
95     Returns:
96         np.ndarray: The resized and padded image/mask with dimensions (target_size, target_size).
97     """
98
99     # Get current dimensions
100    h, w = img.shape[:2]
101
102    # Compute scale factor to fit the longest side
103    scale = target_size / max(w, h)
104    new_w, new_h = int(w * scale), int(h * scale) # New dimensions
105
106    # Select interpolation method (NEAREST for masks, BICUBIC for images)
107    interpolation = cv.INTER_NEAREST if is_mask else cv.INTER_CUBIC
108
109    # Resize image
110    img_resized = cv.resize(img, (new_w, new_h), interpolation=interpolation)
111
112    # Create a blank canvas with padding
113    if len(img.shape) == 3: # RGB image
114        padded_img = np.full((target_size, target_size, 3), fill, dtype=img.dtype)
115    else: # Grayscale/mask
116        padded_img = np.full((target_size, target_size), fill, dtype=img.dtype)
117
118    # Compute padding offsets to center the image
119    paste_x = (target_size - new_w) // 2
120    paste_y = (target_size - new_h) // 2
121
122    # Place the resized image onto the padded canvas
123    padded_img[paste_y:paste_y + new_h, paste_x:paste_x + new_w] = img_resized
124
125    return padded_img
126
127 clip_transform = transforms.Compose([
128     transforms.ToTensor(), # Converts image to [0, 1] range
129     transforms.Normalize(mean=[0.48145466, 0.4578275, 0.40821073],
130         std=[0.26862954, 0.26130258, 0.27577711]) # Converts to ~[-1, 1]
131 ])
132
133 standard_transform = transforms.Compose([
134     transforms.ToTensor()
135 ])
136
137 # -----
138 # Data Augmentation Utilities
139 # -----
140
141 def augmentor(image: np.ndarray, mask: np.ndarray) -> dict:
142     """
143     Apply Albumentations-based augmentations to both image and mask.
144
145     Args:
146         image (np.ndarray): The input image (H, W, C).
147         mask (np.ndarray): The segmentation mask (H, W).
148
149     Returns:
150         dict: Dictionary containing:
151             - 'image' (np.ndarray): The augmented image (H, W, C).
152             - 'mask' (np.ndarray): The augmented mask (H, W).
153     """
154     transform = A.Compose([
155         A.HorizontalFlip(p=0.5), # Flip 50% of the time
156         A.RandomBrightnessContrast(p=0.2), # Adjust brightness & contrast

```



```

157     A.Affine(
158         scale=(0.9, 1.1),
159         translate_percent=(-0.0625, 0.0625),
160         rotate=(-15, 15),
161         interpolation=0, # cv2.INTER_NEAREST (Ensures mask values stay discrete)
162         p=0.5
163     ),
164     A.GaussianBlur(blur_limit=(3, 5), p=0.2), # Blur occasionally
165     A.ColorJitter(brightness=0.2, contrast=0.2, saturation=0.2, hue=0.1, p=0.3) # Color
        jitter
166 ])
167
168 # Apply transformations
169 augmented = transform(image=image, mask=mask)
170
171 return augmented
172
173 #####
174 # Create Prompt-wise Mask #
175 #####
176 def create_point_wise_mask(point_classes, gt_masks):
177     """
178     Create a point-wise mask based on the prompt point and its class.
179
180     Args:
181         prompt_point (torch.Tensor): Tensor of shape (B, 2) with x,y coordinates
182         point_class (torch.Tensor): Tensor of shape (B,) with class indices
183         gt_mask (torch.Tensor): Ground truth mask of shape (B, 1, H, W)
184
185     Returns:
186         torch.Tensor: Binary mask showing regions that should belong to the same class
187                       as the prompt point
188     """
189     B, _, H, W = gt_masks.shape
190     point_wise_masks = []
191
192     for b in range(B):
193         cls = point_classes[b]
194
195         # Create binary mask where 1s represent pixels of the same class as the prompt point
196         point_wise_mask = (gt_masks[b, 0] == cls).float()
197
198         point_wise_masks.append(point_wise_mask)
199
200     return torch.stack(point_wise_masks).unsqueeze(1).long() # (B, 1, H, W)

```

C.3. Preprocessing Runners *preprocess_data.py* and *preprocess_data_with_prompt.py*

These are command-line preprocessing scripts that allow configurable options such as input/output directories, augmentation iterations, image resizing, and prompt point sampling for the prompt-based model.

preprocess_data.py

```

1 # src/preprocess_dataset.py
2
3 from preprocessing import (
4     resize_with_padding,
5     color2class,
6     augmentor
7 )
8 from pathlib import Path
9 from PIL import Image
10 import numpy as np
11 from tqdm import tqdm
12 import torch
13 import argparse
14
15 def prepare_final_dataset(raw_img_dir, raw_msk_dir,

```

```

16         processed_img_dir, processed_msk_dir,
17         base_augmentations, extra_augmentations,
18         target_size):
19     """
20     Runs augmentation and resizing, creating the final training dataset.
21     """
22
23     # Step 0: Initialize the paths of raw images and raw masks into lists
24     raw_img_paths = sorted(Path(raw_img_dir).glob("*."))
25     print(len(raw_img_paths))
26     raw_msk_paths = sorted(Path(raw_msk_dir).glob("*."))
27
28     # create the save directories for processing images and masks
29     processed_img_dir.mkdir(parents=True, exist_ok=True)
30     processed_msk_dir.mkdir(parents=True, exist_ok=True)
31
32     print("Starting_dataset_preparation...")
33     for img_path, msk_path in tqdm(zip(raw_img_paths, raw_msk_paths),
34                                   total=len(raw_msk_paths),
35                                   desc="Final_resizing_&_augmentation",
36                                   unit='image'):
37         # Extract the names of raw image and mask
38         img_name = img_path.stem
39         msk_name = msk_path.stem
40         # Load image and mask
41         img = np.array(Image.open(img_path).convert("RGB"))
42         msk = np.array(Image.open(msk_path).convert("RGB"))
43
44         # Step 1: Resize with padding
45         img_resized = resize_with_padding(img, target_size=target_size, is_mask=False)
46         msk_resized = resize_with_padding(msk, target_size=target_size, is_mask=True)
47
48         msk_class = color2class(msk_resized)
49         unique_classes = torch.unique(torch.tensor(msk_class)).tolist()
50         num_augmentations = extra_augmentations if 1 in unique_classes else base_augmentations
51         for i in range(num_augmentations):
52             augmented = augmentor(img_resized, msk_resized)
53             aug_image, aug_mask = augmented['image'], augmented['mask']
54
55             Image.fromarray(aug_image.astype(np.uint8)).save(
56                 Path(processed_img_dir) / f"aug_{img_name}_{i+1}.jpg")
57
58             Image.fromarray(aug_mask.astype(np.uint8)).save(
59                 Path(processed_msk_dir) / f"aug_{msk_name}_{i+1}.png")
60
61 if __name__ == "__main__":
62     parser = argparse.ArgumentParser(description="Prepare_the_final_dataset_for_segmentation_
63     training.")
64     parser.add_argument('--img_dir', default="./Dataset/Train/color",
65                         help="Directory_containing_raw_images")
66     parser.add_argument('--msk_dir', default="./Dataset/Train/label",
67                         help="Directory_containing_raw_masks")
68     parser.add_argument('--processed_img_dir', default="./Dataset/TrainProcessed/color",
69                         help="Directory_to_save_processed_images_and_points")
70     parser.add_argument('--processed_msk_dir', default="./Dataset/TrainProcessed/label",
71                         help="Directory_to_save_processed_masks")
72     parser.add_argument('--base_augmentations', default=4,
73                         help="Number_of_base_augmentations")
74     parser.add_argument('--extra_augmentations', default=8,
75                         help="Number_of_extra_augmentations")
76     parser.add_argument('--target_size', default=512,
77                         help="Target_size_for_resizing_images_and_masks")
78     args = parser.parse_args()
79     prepare_final_dataset(
80         raw_img_dir=Path(args.img_dir),
81         raw_msk_dir=Path(args.msk_dir),
82         processed_img_dir=Path(args.processed_img_dir),
83         processed_msk_dir=Path(args.processed_msk_dir),

```

```

83     base_augmentations=args.base_augmentations,
84     extra_augmentations=args.extra_augmentations,
85     target_size=args.target_size
86 )

```

preprocess_data_with_prompt.py

```

1  #!/usr/bin/env python
2
3  import numpy as np
4  from pathlib import Path
5  from PIL import Image
6  import torch
7  import argparse
8  from tqdm.contrib.concurrent import process_map # tqdm multiprocessing support
9
10 from data.preprocessing import (
11     resize_with_padding,
12     color2class,
13     augmentor
14 )
15
16 def sample_prompt_type_and_count(max_points=1, box_prob=0.5):
17     """
18     Returns (prompt_type, count)
19     - prompt_type: "points" or "box"
20     - count: number of points (1 for "points", 2 for "box")
21     """
22     if np.random.rand() < box_prob:
23         return "box", 2
24     return "points", 1
25
26
27
28
29 def sample_point_on_mask(mask, num_points, bg_ratio=0.5):
30     mask = color2class(mask)
31     unique_classes = np.unique(mask)
32
33     # ensure the background class (0) is included
34     class_weights = []
35     class_pixel_lists = []
36
37     for cls in unique_classes:
38         class_pixels = np.argwhere(mask == cls)
39         if len(class_pixels) == 0:
40             continue
41         class_pixel_lists.append((cls, class_pixels))
42
43         weight = 1.0 if cls != 0 else bg_ratio
44         class_weights.append(weight * len(class_pixels))
45
46     if len(class_pixel_lists) == 0:
47         return []
48
49     # normalize class weights
50     class_weights = np.array(class_weights, dtype=np.float32)
51     class_weights /= class_weights.sum()
52
53     sampled_points = []
54     for _ in range(num_points):
55         selected_class_idx = np.random.choice(len(class_pixel_lists), p=class_weights)
56         selected_class, class_pixels = class_pixel_lists[selected_class_idx]
57         point_idx = np.random.randint(0, len(class_pixels))
58         y, x = class_pixels[point_idx]
59         sampled_points.append((x, y, selected_class))
60     return sampled_points
61

```

```

62
63
64
65 def process_one_image(args_tuple):
66     (img_path, msk_path, processed_img_dir, processed_msk_dir, processed_point_dir,
67      base_augmentations, extra_augmentations, target_size, points_per_image, mode, box_prob) =
68         args_tuple
69
70     img_name = img_path.stem
71     msk_name = msk_path.stem
72
73     img = np.array(Image.open(img_path).convert("RGB"))
74     msk = np.array(Image.open(msk_path).convert("RGB"))
75
76     img_resized = resize_with_padding(img, target_size=target_size, is_mask=False)
77     msk_resized = resize_with_padding(msk, target_size=target_size, is_mask=True)
78
79     if mode in [1, 2, 3]: # test modes: no augmentation
80         Image.fromarray(img_resized.astype(np.uint8)).save(processed_img_dir / f"{img_name}.jpg")
81         Image.fromarray(msk_resized.astype(np.uint8)).save(processed_msk_dir / f"{msk_name}.png")
82
83     prompt_type, num = sample_prompt_type_and_count(points_per_image, box_prob=box_prob)
84     # Here we sample points and each point is a tuple (x, y, cls)
85     sampled_points = sample_point_on_mask(msk_resized, num_points=num, bg_ratio=0.2)
86
87     points_file_path = processed_point_dir / f"{img_name}_points.txt"
88     with open(points_file_path, 'w') as f:
89         if prompt_type == "box":
90             # Compute box using first two sampled points
91             x_coords = [pt[0] for pt in sampled_points[:2]]
92             y_coords = [pt[1] for pt in sampled_points[:2]]
93             x1, x2 = sorted(x_coords)
94             y1, y2 = sorted(y_coords)
95             # Convert box region to class labels before processing
96             box_mask_class = color2class(msk_resized[y1:y2+1, x1:x2+1])
97             if box_mask_class.size == 0:
98                 box_cls = -1
99             else:
100                 cls_counts = np.bincount(box_mask_class.flatten())
101                 box_cls = np.argmax(cls_counts) if len(cls_counts) > 0 else -1
102
103             with open(points_file_path, 'w') as f:
104                 f.write(f"{x1},{y1},{x2},{y2},{box_cls}\n")
105         else:
106             # For single point prompt
107             x, y, cls = sampled_points[0]
108             f.write(f"{x},{y},{cls}\n")
109
110     return
111
112 # mode == 0: train with augmentation
113 msk_class = color2class(msk_resized)
114 unique_classes = torch.unique(torch.tensor(msk_class)).tolist()
115 num_augmentations = extra_augmentations if 1 in unique_classes else base_augmentations
116
117 for i in range(num_augmentations):
118     augmented = augmentor(img_resized, msk_resized)
119     aug_image, aug_mask = augmented['image'], augmented['mask']
120
121     prompt_type, num = sample_prompt_type_and_count(points_per_image, box_prob=box_prob)
122     sampled_points = sample_point_on_mask(aug_mask, num_points=num)
123
124     aug_img_path = processed_img_dir / f"aug_{img_name}_{i+1}.jpg"
125     aug_msk_path = processed_msk_dir / f"aug_{msk_name}_{i+1}.png"
126     points_file_path = processed_point_dir / f"aug_{img_name}_{i+1}_points.txt"
127
128     Image.fromarray(aug_image.astype(np.uint8)).save(aug_img_path)
129     Image.fromarray(aug_mask.astype(np.uint8)).save(aug_msk_path)

```

```

129     with open(points_file_path, 'w') as f:
130         if prompt_type == "box":
131             # Compute box using first two sampled points
132             x_coords = [pt[0] for pt in sampled_points[:2]]
133             y_coords = [pt[1] for pt in sampled_points[:2]]
134             x1, x2 = sorted(x_coords)
135             y1, y2 = sorted(y_coords)
136             # Extract the mask region to compute the mode class
137             box_mask_class = color2class(aug_mask[y1:y2+1, x1:x2+1])
138             if box_mask_class.size == 0:
139                 box_cls = -1
140             else:
141                 box_cls_counts = np.bincount(box_mask_class.flatten())
142                 if len(box_cls_counts) == 0:
143                     box_cls = -1
144                 else:
145                     box_cls = np.argmax(box_cls_counts)
146             f.write(f"{x1},{y1},{x2},{y2},{box_cls}\n")
147         else:
148             # For single point prompt
149             x, y, cls = sampled_points[0]
150             f.write(f"{x},{y},{cls}\n")
151
152
153
154
155 def prepare_final_dataset(raw_img_dir, raw_msk_dir,
156                          processed_img_dir, processed_msk_dir,
157                          base_augmentations, extra_augmentations,
158                          target_size, points_per_image=5, mode=0, num_workers=4, box_prob=0.2):
159     raw_img_paths = sorted(Path(raw_img_dir).glob("*."))
160     raw_msk_paths = sorted(Path(raw_msk_dir).glob("*."))
161
162     processed_img_dir = Path(processed_img_dir)
163     processed_msk_dir = Path(processed_msk_dir)
164     processed_point_dir = processed_img_dir / "points"
165     processed_img_dir.mkdir(parents=True, exist_ok=True)
166     processed_msk_dir.mkdir(parents=True, exist_ok=True)
167     processed_point_dir.mkdir(parents=True, exist_ok=True)
168
169     print(f"Preparing_{len(raw_img_paths)}_image-mask_pairs_using_{num_workers}_workers...")
170
171     args_list = [
172         (img_path, msk_path, processed_img_dir, processed_msk_dir, processed_point_dir,
173          base_augmentations, extra_augmentations, target_size, points_per_image, mode, box_prob)
174         for img_path, msk_path in zip(raw_img_paths, raw_msk_paths)
175     ]
176
177
178     process_map(process_one_image, args_list, max_workers=num_workers, chunksize=1, desc="
179                 Processing")
180
181 if __name__ == "__main__":
182     parser = argparse.ArgumentParser()
183     parser.add_argument('--mode', type=int, default=0, choices=[0, 1, 2, 3],
184                         help="0=_full_preprocessing_with_augmentation, 1=_prompt_point_only, 2=_test_set_(no_
185                             augmentation, _with_class), 3=_test_set_(no_augmentation, _no_class)")
186     parser.add_argument('--img_dir', default="./Dataset/TrainVal/color",
187                         help="Directory_containing_raw_images")
188     parser.add_argument('--msk_dir', default="./Dataset/TrainVal/label",
189                         help="Directory_containing_raw_masks")
190     parser.add_argument('--processed_img_dir', default="./Dataset/ProcessedWithPrompt_test/color"
191                         ,
192                         help="Directory_to_save_processed_images_and_points")
193     parser.add_argument('--processed_msk_dir', default="./Dataset/ProcessedWithPrompt_test/label"
194                         ,
195                         help="Directory_to_save_processed_masks")
196     parser.add_argument('--points_per_image', type=int, default=1,

```



```

193     help="Number_of_prompt_points_to_generate")
194     parser.add_argument('--num_workers', type=int, default=8,
195         help="Number_of_parallel_workers_for_processing")
196     parser.add_argument('--box_prob', type=float, default=0.5,
197         help="Probability_of_sampling_a_box_instead_of_points")
198
199     args = parser.parse_args()
200
201     prepare_final_dataset(
202         raw_img_dir=Path(args.img_dir),
203         raw_msk_dir=Path(args.msk_dir),
204         processed_img_dir=Path(args.processed_img_dir),
205         processed_msk_dir=Path(args.processed_msk_dir),
206         base_augmentations=4,
207         extra_augmentations=8,
208         target_size=512,
209         points_per_image=args.points_per_image,
210         mode=args.mode,
211         num_workers=args.num_workers,
212         box_prob=args.box_prob
213     )

```

C.4. Train/Validation Split Utility *train_val_split.py*

This script partitions the dataset into training and validation subsets. The test set remains completely unseen to ensure fair evaluation, and the best model is selected based on performance on the validation set.

train_val_split.py

```

1  from pathlib import Path
2
3  source_color_dir = Path("Dataset/TrainVal/color/")
4  source_label_dir = Path("Dataset/TrainVal/label/")
5
6  import random
7  from pathlib import Path
8  import shutil
9
10 # Set seed for reproducibility
11 random.seed(42)
12
13 # Source directories
14 source_color_dir = Path("Dataset/TrainVal/color/")
15 source_label_dir = Path("Dataset/TrainVal/label/")
16
17 # Destination directories
18 train_color_dir = Path("Dataset/Train/color/")
19 train_label_dir = Path("Dataset/Train/label/")
20 val_color_dir = Path("Dataset/Val/color/")
21 val_label_dir = Path("Dataset/Val/label/")
22
23 # Create destination folders
24 for path in [train_color_dir, train_label_dir, val_color_dir, val_label_dir]:
25     path.mkdir(parents=True, exist_ok=True)
26
27 # Get list of color image files (assuming same names for labels)
28 color_files = sorted(source_color_dir.glob("*"))
29 total = len(color_files)
30 split_idx = int(total * 0.9)
31
32 # Shuffle and split
33 random.shuffle(color_files)
34 train_files = color_files[:split_idx]
35 val_files = color_files[split_idx:]
36
37 # Helper to copy files
38 def copy_files(file_list, color_dest, label_dest):
39     for color_file in file_list:

```

```

40     label_filename = color_file.stem + ".png"
41     label_file = source_label_dir / label_filename
42     shutil.copy2(color_file, color_dest / color_file.name)
43     shutil.copy2(label_file, label_dest / label_file.name)
44
45 # Copy files to train and val directories
46 copy_files(train_files, train_color_dir, train_label_dir)
47 copy_files(val_files, val_color_dir, val_label_dir)
48
49 print(f"Total_images:_{total}")
50 print(f"Training_images:_{len(train_files)}")
51 print(f"Validation_images:_{len(val_files)}")

```

D. Model Structures

This appendix provides the implementation details of the four segmentation models compared in this report. These model definitions serve as supporting material for Section 3, where the design rationale and high-level architecture are discussed.

D.1. U-Net Model

unet_segmentation.py

```

1 import torch
2 import torch.nn as nn
3
4 class UNet(nn.Module):
5     def __init__(self, in_channels=3, out_channels=3, init_features=64):
6         super(UNet, self).__init__()
7
8         features = init_features
9         # Encoder blocks
10        self.encoder1 = self._conv_block(in_channels, features)
11        self.encoder2 = self._conv_block(features, features * 2)
12        self.encoder3 = self._conv_block(features * 2, features * 4)
13        self.encoder4 = self._conv_block(features * 4, features * 8)
14
15        # Bottleneck
16        self.bottleneck = self._conv_block(features * 8, features * 16)
17
18        # Decoder upsampling layers (now ConvTranspose2d)
19        self.up4 = nn.ConvTranspose2d(features * 16, features * 8, kernel_size=2, stride=2)
20        self.decoder4 = self._conv_block(features * 16, features * 8)
21
22        self.up3 = nn.ConvTranspose2d(features * 8, features * 4, kernel_size=2, stride=2)
23        self.decoder3 = self._conv_block(features * 8, features * 4)
24
25        self.up2 = nn.ConvTranspose2d(features * 4, features * 2, kernel_size=2, stride=2)
26        self.decoder2 = self._conv_block(features * 4, features * 2)
27
28        self.up1 = nn.ConvTranspose2d(features * 2, features, kernel_size=2, stride=2)
29        self.decoder1 = self._conv_block(features * 2, features)
30
31        # Final output layer
32        self.final_layer = nn.Conv2d(features, out_channels, kernel_size=1)
33
34        # Max pooling
35        self.pool = nn.MaxPool2d(kernel_size=2, stride=2)
36
37    def _conv_block(self, in_channels, out_channels):
38        return nn.Sequential(
39            nn.Conv2d(in_channels, out_channels, kernel_size=3, padding=1),
40            nn.BatchNorm2d(out_channels),
41            nn.ReLU(inplace=True),
42            nn.Conv2d(out_channels, out_channels, kernel_size=3, padding=1),
43            nn.BatchNorm2d(out_channels),
44            nn.ReLU(inplace=True),

```

```

45         )
46
47     def forward(self, x):
48         # Encoder
49         enc1 = self.encoder1(x) # (B, features, 512, 512)
50         enc2 = self.encoder2(self.pool(enc1)) # (B, features*2, 256, 256)
51         enc3 = self.encoder3(self.pool(enc2)) # (B, features*4, 128, 128)
52         enc4 = self.encoder4(self.pool(enc3)) # (B, features*8, 64, 64)
53
54         # Bottleneck
55         bottleneck = self.bottleneck(self.pool(enc4)) # (B, features*16, 32, 32)
56
57         # Decoder with skip connections
58         dec4 = self.up4(bottleneck) # (B, features*8, 64, 64)
59         dec4 = torch.cat((dec4, enc4), dim=1) # (B, features*16, 64, 64)
60         dec4 = self.decoder4(dec4)
61
62         dec3 = self.up3(dec4) # (B, features*4, 128, 128)
63         dec3 = torch.cat((dec3, enc3), dim=1) # (B, features*8, 128, 128)
64         dec3 = self.decoder3(dec3)
65
66         dec2 = self.up2(dec3) # (B, features*2, 256, 256)
67         dec2 = torch.cat((dec2, enc2), dim=1) # (B, features*4, 256, 256)
68         dec2 = self.decoder2(dec2)
69
70         dec1 = self.up1(dec2) # (B, features, 512, 512)
71         dec1 = torch.cat((dec1, enc1), dim=1) # (B, features*2, 512, 512)
72         dec1 = self.decoder1(dec1)
73
74         # Final output layer
75         return self.final_layer(dec1)

```

D.2. Autoencoder-based Model

autoencoder_segmentation.py

```

1  import torch
2  import torch.nn as nn
3  import torch.optim as optim
4  from pathlib import Path
5  from tqdm import tqdm
6
7  #####
8  ## Autoencoder Architecture ##
9  #####
10
11 class AutoEncoder(nn.Module):
12     """
13     Autoencoder architecture for image reconstruction.
14     The encoder compresses the input image into a latent representation,
15     and the decoder reconstructs the image from this representation.
16     Args:
17         init_features (int): Number of initial features for the encoder.
18     """
19     def __init__(self, init_features=128):
20         super(AutoEncoder, self).__init__()
21
22         f = init_features
23
24         # Encoder
25         self.encoder = nn.Sequential(
26             nn.Conv2d(3, f, kernel_size=3, stride=2, padding=1), # 3*512*512 > 64*256*256
27             nn.BatchNorm2d(f),
28             nn.ReLU(),
29             nn.Conv2d(f, f * 2, kernel_size=3, stride=2, padding=1), # 64*256*256 > 128*128*128
30             nn.BatchNorm2d(f * 2),
31             nn.ReLU(),

```

```

32         nn.Conv2d(f * 2, f * 4, kernel_size=3, stride=2, padding=1), # 128*128*128 >
           256*64*64
33         nn.BatchNorm2d(f * 4),
34         nn.ReLU(),
35         nn.Conv2d(f * 4, f * 8, kernel_size=3, stride=2, padding=1), # 256*64*64 > 512*32*32
36         nn.BatchNorm2d(f * 8),
37         nn.ReLU(),
38         nn.Conv2d(f * 8, f * 16, kernel_size=3, stride=2, padding=1), #512*32*32 > 1024*16*16
39     )
40
41     # Decoder
42     self.decoder = nn.Sequential(
43         nn.ConvTranspose2d(f * 16, f * 8, kernel_size=3, stride=2, padding=1, output_padding
           =1), # 1024*16*16 > 512*32*32
44         nn.BatchNorm2d(f * 8),
45         nn.ReLU(),
46         nn.ConvTranspose2d(f * 8, f * 4, kernel_size=3, stride=2, padding=1, output_padding
           =1), # 512*32*32 > 256*64*64
47         nn.BatchNorm2d(f * 4),
48         nn.ReLU(),
49         nn.ConvTranspose2d(f * 4, f * 2, kernel_size=3, stride=2, padding=1, output_padding
           =1), # 256*64*64 > 128*128*128
50         nn.BatchNorm2d(f * 2),
51         nn.ReLU(),
52         nn.ConvTranspose2d(f * 2, f, kernel_size=3, stride=2, padding=1, output_padding=1),
           # 128*128*128 > 64*256*256
53         nn.BatchNorm2d(f),
54         nn.ReLU(),
55         nn.ConvTranspose2d(f, 3, kernel_size=3, stride=2, padding=1, output_padding=1),
           # 64*256*256 > 3*512*512
56         nn.Sigmoid()
57     )
58
59     def forward(self, x):
60         latent = self.encoder(x)
61         reconstructed = self.decoder(latent)
62         return reconstructed
63
64     #####
65     ## Segmentation Decoder Head ##
66     #####
67
68     class SegmentationDecoder(nn.Module):
69         '''
70         Segmentation Decoder for the Autoencoder.
71         Args:
72             input_channel (int): Number of input channels from the encoder.
73             output_channel (int): Number of output channels for segmentation.
74         '''
75         def __init__(self, input_channel, output_channel):
76             super().__init__()
77             self.seg_decoder = nn.Sequential(
78                 nn.ConvTranspose2d(input_channel, 512, 3, stride=2, padding=1, output_padding=1), #
           1024*16*16 > 512*32*32
79                 nn.BatchNorm2d(512),
80                 nn.ReLU(),
81                 nn.ConvTranspose2d(512, 256, 3, stride=2, padding=1, output_padding=1), # 512*32*32
           > 256*64*64
82                 nn.BatchNorm2d(256),
83                 nn.ReLU(),
84                 nn.ConvTranspose2d(256, 128, 3, stride=2, padding=1, output_padding=1), # 256*64*64
           > 128*128*128
85                 nn.BatchNorm2d(128),
86                 nn.ReLU(),
87                 nn.ConvTranspose2d(128, 64, 3, stride=2, padding=1, output_padding=1), #
           128*128*128 > 64*256*256
88                 nn.BatchNorm2d(64),
89                 nn.ReLU(),

```

```

90         nn.ConvTranspose2d(64, output_channel, 3, stride=2, padding=1, output_padding=1), #
          64*256*256 > 3*512*512
91     )
92
93     def forward(self, features):
94         return self.seg_decoder(features)
95
96     #####
97     ## Final Autoencoder-based Segmentation ##
98     #####
99
100 class AutoEncoderSegmentation(nn.Module):
101     '''
102     Autoencoder-based Segmentation Model.
103     Args:
104         encoder (nn.Module): Encoder part of the autoencoder.
105         num_classes (int): Number of classes for segmentation.
106     Returns:
107         out (torch.Tensor): Segmentation output.
108     '''
109     def __init__(self, encoder, num_classes=3):
110         super().__init__()
111         # Frozen encoder from pre-trained autoencoder
112         self.encoder = encoder
113         for param in self.encoder.parameters():
114             param.requires_grad = False
115         # Segmentation Decoder
116         self.seg_decoder = SegmentationDecoder(input_channel=self.encoder[-1].out_channels,
117                                                output_channel=num_classes)
118
119     def forward(self, x):
120         with torch.no_grad():
121             features = self.encoder(x) # Extract latent features
122             out = self.seg_decoder(features)
123             return out
124
125     #####
126     ## Autoencoder Pretraining Function ##
127     #####
128     def pretrain_autoencoder(autoencoder, train_loader, val_loader,
129                             save_dir, save_name,
130                             num_epochs, device, patience):
131         '''
132         Pretrain the autoencoder on the training dataset.
133         Args:
134             autoencoder (nn.Module): Autoencoder model.
135             train_loader (DataLoader): DataLoader for training data.
136             val_loader (DataLoader): DataLoader for validation data.
137             save_dir (str): Directory to save the pretrained model.
138             save_name (str): Name of the saved model file.
139             num_epochs (int): Number of epochs for pretraining.
140             device (str): Device to use ("cuda" or "cpu").
141             patience (int): Number of epochs with no improvement after which training will be stopped
142
143         Returns:
144             history (dict): Dictionary containing training and validation losses.
145         '''
146         autoencoder.to(device)
147         criterion = nn.MSELoss()
148         optimizer = optim.Adam(autoencoder.parameters(), lr=1e-3)
149         best_val_loss = float('inf')
150         best_model_state = None
151         epochs_no_improve = 0
152         save_path = Path(save_dir) / save_name
153
154         # Initialize history dict to store losses
155         history = {"train_loss": [], "val_loss": []}

```

```

155
156
157 for epoch in range(num_epochs):
158     # Training Phase
159     autoencoder.train()
160     running_loss = 0.0
161     for images, _, _ in tqdm(train_loader, desc=f"Pretrain_Epoch_{epoch+1}/{num_epochs}"):
162         images = images.to(device)
163         optimizer.zero_grad()
164         reconstructed = autoencoder(images)
165         loss = criterion(reconstructed, images)
166         loss.backward()
167         optimizer.step()
168         running_loss += loss.item()
169
170     avg_train_loss = running_loss / len(train_loader)
171
172     # Validation Phase
173     autoencoder.eval()
174     val_running_loss = 0.0
175     with torch.no_grad():
176         for images, _, _ in tqdm(val_loader, desc=f"Pretrain_Epoch_{epoch+1}/{num_epochs}_
177             [Val]", unit='image'):
178             images = images.to(device)
179             reconstructed = autoencoder(images)
180             val_loss = criterion(reconstructed, images)
181             val_running_loss += val_loss.item()
182
183     avg_val_loss = val_running_loss / len(val_loader)
184
185     # Store the losses for plotting
186     history["train_loss"].append(avg_train_loss)
187     history["val_loss"].append(avg_val_loss)
188     # Print Losses
189     print(f"Epoch_{epoch+1}/{num_epochs}_|_Train_Loss:_{avg_train_loss:.6f}_|_Val_Loss:_{
190         avg_val_loss:.6f}")
191
192     # Early Stopping Check (on validation loss)
193     if avg_val_loss < best_val_loss:
194         best_val_loss = avg_val_loss
195         best_model_state = autoencoder.state_dict()
196         epochs_no_improve = 0
197     else:
198         epochs_no_improve += 1
199         print(f"No_improvement_for_{epochs_no_improve}_epoch(s)")
200
201     if epochs_no_improve >= patience:
202         print(f"Early_stopping_triggered_at_epoch_{epoch+1}")
203         print(f"Minimum_Val_MSE_Loss:_{best_val_loss:.6f}")
204         break
205
206     # save the model_state with the best performance
207     torch.save(best_model_state, save_path)
208     print(f"Best_autoencoder_model_saved_at_{save_path}")
209
210 return history

```

D.3. CLIP-featured Model

clip_segmentation.py

```

1 # The related blocks of CLIP segmentation models are store all-in-one here
2
3 import torch
4 import torch.nn as nn
5 #####
6 # Feature Extraction #
7 #####
8

```

```

9 class CLIPFeatureExtractor(nn.Module):
10     def __init__(self, clip_model):
11         super().__init__()
12         visual = clip_model.visual
13         self.conv1 = visual.conv1
14         self.bn1 = visual.bn1
15         self.relu1 = visual.relu1
16         self.conv2 = visual.conv2
17         self.bn2 = visual.bn2
18         self.relu2 = visual.relu2
19         self.conv3 = visual.conv3
20         self.bn3 = visual.bn3
21         self.relu3 = visual.relu3
22
23         self.layer1 = visual.layer1
24         self.layer2 = visual.layer2
25         self.layer3 = visual.layer3
26         self.layer4 = visual.layer4
27
28     def forward(self, x):
29         x = self.conv1(x)
30         x = self.bn1(x)
31         x = self.relu1(x)
32         x = self.conv2(x)
33         x = self.bn2(x)
34         x = self.relu2(x)
35         x = self.conv3(x)
36         x = self.bn3(x)
37         x = self.relu3(x)
38         x = self.layer1(x)
39         x = self.layer2(x)
40         x = self.layer3(x)
41         x = self.layer4(x)
42         return x
43
44 #####
45 ## Segmentation Head ##
46 #####
47
48 class CLIPSegmentationHead(nn.Module):
49     """
50     Updated segmentation head for ModifiedResNet in CLIP.
51     Adjusted for 2048 input channels and 14x14 input spatial size.
52     """
53     def __init__(self, input_channels, output_channels):
54         super().__init__()
55
56         # Reduce channels first
57         self.reduce_channels = nn.Sequential(
58             nn.Conv2d(input_channels, 1024, kernel_size=1), # Reduce channels 2048 > 1024
59             nn.BatchNorm2d(1024),
60             nn.ReLU()
61         )
62
63         # Upsampling layers
64         self.upsample1 = nn.Sequential(
65             nn.ConvTranspose2d(1024, 512, kernel_size=3, stride=2, padding=1, output_padding=1),
66             # 14x14 > 28x28
67             nn.BatchNorm2d(512),
68             nn.ReLU()
69         )
70
71         self.upsample2 = nn.Sequential(
72             nn.Upsample(scale_factor=8, mode='bilinear', align_corners=True), # 112x112 > 224
73             # x224
74             nn.Conv2d(512, 64, kernel_size=3, padding=1),
75             nn.BatchNorm2d(64),
76             nn.ReLU()

```

```

75     )
76
77     # Final segmentation layer
78     self.final_conv = nn.Conv2d(64, output_channels, kernel_size=1) # Output segmentation
79     map
80
81     def forward(self, x):
82         x = self.reduce_channels(x) # Reduce channels: (B, 20248, 14, 14) > (B, 1024, 14, 14)
83         x1 = self.upsample1(x) # 14x14 > 28x28
84         x2 = self.upsample2(x1) # 112x112 > 224x224
85         out = self.final_conv(x2)
86
87         return out
88
89 #####
90 ### Final CLIP Segmentation model ###
91 #####
92 class CLIPSegmentationModel(nn.Module):
93     """
94     Lightweight CLIP-based segmentation model using only the necessary vision layers from RN50.
95     """
96     def __init__(self, clip_model, num_classes):
97         super().__init__()
98
99         # Replace full CLIP model with a custom extractor that only includes needed layers
100         self.feature_extractor = CLIPFeatureExtractor(clip_model) # defined above
101         for param in self.feature_extractor.parameters():
102             param.requires_grad = False # Freeze backbone
103
104         self.seg_head = CLIPSegmentationHead(input_channels=2048, output_channels=num_classes)
105
106     def forward(self, x):
107         with torch.no_grad():
108             x = self.feature_extractor(x)
109             x = self.seg_head(x)
110         return x

```

D.4. Prompt-based Model

prompt_segmentation.py

```

1 import torch
2 import torch.nn as nn
3 from models.unet_segmentation import UNet
4
5 class UNetEncoder(nn.Module):
6     def __init__(self, in_channels, init_features):
7         super().__init__()
8         f = init_features
9         self.encoder1 = self._conv_block(in_channels, f)
10        self.encoder2 = self._conv_block(f, f * 2)
11        self.encoder3 = self._conv_block(f * 2, f * 4)
12        self.encoder4 = self._conv_block(f * 4, f * 8)
13        self.bottleneck = self._conv_block(f * 8, f * 16)
14        self.bottleneck_channels = f * 16
15        self.pool = nn.MaxPool2d(kernel_size=2, stride=2)
16
17    def _conv_block(self, in_c, out_c):
18        return nn.Sequential(
19            nn.Conv2d(in_c, out_c, 3, padding=1),
20            nn.BatchNorm2d(out_c),
21            nn.ReLU(inplace=True),
22            nn.Conv2d(out_c, out_c, 3, padding=1),
23            nn.BatchNorm2d(out_c),
24            nn.ReLU(inplace=True),
25        )
26

```



```

27     def forward(self, x):
28         enc1 = self.encoder1(x)
29         enc2 = self.encoder2(self.pool(enc1))
30         enc3 = self.encoder3(self.pool(enc2))
31         enc4 = self.encoder4(self.pool(enc3))
32         bottleneck = self.bottleneck(self.pool(enc4))
33         return [enc1, enc2, enc3, enc4], bottleneck
34
35 class UNetDecoder(nn.Module):
36     def __init__(self, bottleneck_channels, out_channels, init_features):
37         super().__init__()
38         f = init_features
39
40         self.up4 = nn.ConvTranspose2d(bottleneck_channels, f * 8, kernel_size=2, stride=2)
41         self.dec4 = self._conv_block(f * 8 + f * 8, f * 8)
42
43         self.up3 = nn.ConvTranspose2d(f * 8, f * 4, kernel_size=2, stride=2)
44         self.dec3 = self._conv_block(f * 4 + f * 4, f * 4)
45
46         self.up2 = nn.ConvTranspose2d(f * 4, f * 2, kernel_size=2, stride=2)
47         self.dec2 = self._conv_block(f * 2 + f * 2, f * 2)
48
49         self.up1 = nn.ConvTranspose2d(f * 2, f, kernel_size=2, stride=2)
50         self.dec1 = self._conv_block(f + f, f)
51
52         self.final = nn.Conv2d(f, out_channels, kernel_size=1)
53
54     def _conv_block(self, in_c, out_c):
55         return nn.Sequential(
56             nn.Conv2d(in_c, out_c, 3, padding=1),
57             nn.BatchNorm2d(out_c),
58             nn.ReLU(inplace=True),
59             nn.Conv2d(out_c, out_c, 3, padding=1),
60             nn.BatchNorm2d(out_c),
61             nn.ReLU(inplace=True),
62         )
63
64     def forward(self, enc_feats, bottleneck):
65         enc1, enc2, enc3, enc4 = enc_feats
66
67         d4 = self.up4(bottleneck)
68         d4 = self.dec4(torch.cat([d4, enc4], dim=1))
69
70         d3 = self.up3(d4)
71         d3 = self.dec3(torch.cat([d3, enc3], dim=1))
72
73         d2 = self.up2(d3)
74         d2 = self.dec2(torch.cat([d2, enc2], dim=1))
75
76         d1 = self.up1(d2)
77         d1 = self.dec1(torch.cat([d1, enc1], dim=1))
78
79         return self.final(d1)
80
81 class PointPromptEncoder(nn.Module):
82     """
83     Encode prompt points into a spatial representation using pre-generated heatmaps
84     """
85     def __init__(self, input_channels=1, output_dim=1024, num_classes=3):
86         """
87         Args:
88             input_channels (int): Number of input channels for the prompt heatmap.
89             output_dim (int): Dimension of the final output representation.
90             num_classes (int): Number of classes for point classification.
91         """
92         super().__init__()
93         # define the intermediate hidden dimension

```

```

95     self.hidden_dim = output_dim // 2
96     # Point class embedding (optional)
97     self.class_embedding = nn.Embedding(num_embeddings=num_classes, embedding_dim=output_dim)
98         # 3 classes (0=background, 1=cat, 2=dog)
99
100     # Spatial encoding mechanism
101     self.spatial_encoder = nn.Sequential(
102         nn.Conv2d(input_channels, self.hidden_dim, kernel_size=3, padding=1, stride=4),
103         nn.ReLU(),
104         nn.Conv2d(self.hidden_dim, output_dim, kernel_size=3, padding=1, stride=4)
105     )
106
107     def forward(self, prompt_heatmap, point_class=None):
108         """
109         Args:
110             prompt_heatmap (torch.Tensor): Pre-generated heatmap from dataset (B, 1, H, W)
111             point_class (torch.Tensor, optional): Class labels for the prompt points (B,)
112
113         Returns:
114             torch.Tensor: Encoded prompt embedding (B, output_dim, H, W)
115         """
116         B = prompt_heatmap.shape[0]
117
118         # Spatially encode heatmap
119         spatial_prompt = self.spatial_encoder(prompt_heatmap)
120
121         # Incorporate class information if provided
122         if point_class is not None:
123             class_embed = self.class_embedding(point_class) # (B, output_dim)
124             class_embed = class_embed.view(B, -1, 1, 1) # (B, output_dim, 1, 1)
125
126             # Broadcast class embedding across spatial dimensions so the final output is globally
127             # class-aware
128             class_embed = class_embed.expand(-1, -1, spatial_prompt.size(2), spatial_prompt.size(3))
129
130             # Combine class information with spatial prompt
131             return spatial_prompt + class_embed
132
133         return spatial_prompt
134
135 class PromptSegmentation(nn.Module):
136     def __init__(self, unet_in_channels=3, prompt_dim=1024, unet_init_features=64):
137         super().__init__()
138
139         # bottleneck features = init_features * 16 -> (64 * 16 = 1024)
140         self.encoder = UNetEncoder(in_channels=unet_in_channels, init_features=unet_init_features)
141
142         # prompt_encoder output channels = 64 * 16 = 1024
143         self.prompt_encoder = PointPromptEncoder(input_channels=1, output_dim=prompt_dim,
144             num_classes=3) # 3 classes (0=background, 1=cat, 2=dog)
145
146         # decoder expects (B, 1024 32, 32) from encoder concatenated (B, 1024, 32, 32) from
147         # prompt_encoder
148         self.decoder = UNetDecoder(bottleneck_channels=self.encoder.bottleneck_channels +
149             prompt_dim,
150                                     out_channels=1, init_features=unet_init_features)
151
152     def forward(self, image, prompt_heatmap, point_class=None):
153         # Step 1: Encode image -> get encoder features and bottleneck
154         enc_features, image_bottleneck = self.encoder(image)
155
156         # Step 2: Encode prompt heatmap
157         prompt_features = self.prompt_encoder(prompt_heatmap, point_class) # (B, prompt_dim, H
158             /16, W/16)
159
160         # Step 3: Concatenate prompt with image bottleneck
161         combined_bottleneck = torch.cat([image_bottleneck, prompt_features], dim=1)

```

```

155
156     # Step 4: Decode
157     mask_logits = self.decoder(enc_features, combined_bottleneck)
158
159     return mask_logits

```

E. Training Utilities

This appendix presents the core components related to model training. It includes modular training functions tailored to the nuances of each model type, visualization tools for tracking training progress, and customized loss functions for robust segmentation performance.

E.1. Training Functions *training.py*

training.py

Provides general training loops, including support for both standard and prompt-based segmentation models. Implements early stopping, optimizer setup, and performance evaluation during training.

```

1 import os
2 from pathlib import Path
3 from tqdm import tqdm
4 import torch
5 from data.preprocessing import create_point_wise_mask
6
7 #####
8 # Segmentaiton Training #
9 #####
10 def training(
11     model, train_loader, val_loader,
12     train_criterion, val_criterion, # Use different loss functions
13     optimizer, num_epochs=10, device="cuda",
14     save_dir=" ../params", save_name=None,
15     patience=5):
16     # assign model
17     model.to(device).float()
18
19     # print model-related info
20     print("    _Model_Structure:\n", model)
21
22     # Count total parameters and trainable parameters
23     total_params = sum(p.numel() for p in model.parameters())
24     trainable_params = sum(p.numel() for p in model.parameters() if p.requires_grad)
25
26     print(f"    _Total_Parameters:_{total_params:,}")
27     print(f"    _Trainable_Parameters:_{trainable_params:,}")
28
29     # Estimate model size (in MB)
30     param_size_MB = total_params * 4 / (1024 ** 2) # float32 => 4 bytes
31     print(f"    _Estimated_Parameter_Size:_{param_size_MB:.2f}_MB")
32
33     # training utils
34     history = {"train_loss": [], "val_loss": []}
35
36     os.makedirs(save_dir, exist_ok=True) # Create checkpoint directory
37     save_path = Path(save_dir) / save_name
38
39     # Initialize Early Stopping Variables
40     best_val_loss = float("inf") # Set to a large value
41     best_model_config = None # To store the best model parameters
42     epochs_no_improve = 0 # Counter for non-improving epochs
43
44     for epoch in range(num_epochs):
45         # Training Phase
46         model.train()
47         running_loss = 0.0
48

```

```

49     for images, masks, _, _ in tqdm(train_loader, desc=f"Epoch_{epoch+1}/{num_epochs}_[Train]"
50                                     , unit='_Batches'):
51         images, masks = images.to(device), masks.to(device)
52
53         optimizer.zero_grad()
54         outputs = model(images)
55         loss = train_criterion(outputs, masks)
56
57         loss.backward()
58         optimizer.step()
59
60         running_loss += loss.item()
61
62     train_loss = running_loss / len(train_loader)
63     history["train_loss"].append(train_loss)
64
65     # Validation Phase
66     model.eval()
67     val_loss = 0.0
68
69     with torch.no_grad():
70         for images, masks, _, _ in tqdm(val_loader, desc=f"Epoch_{epoch+1}/{num_epochs}_[Val]"
71                                         , unit='_Batches'):
72             images, masks = images.to(device), masks.to(device)
73
74             outputs = model(images)
75             loss = val_criterion(outputs, masks)
76
77             val_loss += loss.item()
78
79     val_loss /= len(val_loader)
80     history["val_loss"].append(val_loss)
81     print(f"Epoch_{epoch+1}/{num_epochs}_|_Train_Loss:_{train_loss:.6f}_|_Val_Loss:_{val_loss:.6f}")
82
83     # Early Stopping Logic
84     if val_loss < best_val_loss:
85         best_val_loss = val_loss # Update best loss
86         best_epoch = epoch + 1
87         epochs_no_improve = 0 # Reset counter
88         best_model_config = model.state_dict()
89     else:
90         epochs_no_improve += 1 # Increment counter
91         print(f"No_improvement_for_{epochs_no_improve}_epochs.")
92
93     if epochs_no_improve >= patience:
94         # Save model
95         torch.save(best_model_config, save_path) # Save best model
96         print(f"Early_stopping_triggered_after_{epoch+1}_epochs!")
97         print(f"Best_val_loss_at_epoch_{best_epoch}:_{best_val_loss:.8}")
98         print(f" Best _model_config_saved_at:_{save_path}")
99         return history
100
101     torch.save(best_model_config, save_path) # save the model if the early stop is not triggered
102     print(f" Model _config_saved_at:_{save_path}_in_epoch_{best_epoch}")
103     return history
104
105 def prompt_training(
106     model, train_loader, val_loader,
107     train_criterion, val_criterion,
108     optimizer, num_epochs, device="cuda",
109     save_dir=" ../params", save_name= None,
110     patience=5):
111     model.to(device).float()
112
113     #####
114     # Print model-related info #
115     #####

```

```

114 print("    _Model_Structure:\n", model)
115
116 total_params = sum(p.numel() for p in model.parameters())
117 trainable_params = sum(p.numel() for p in model.parameters() if p.requires_grad)
118 param_size_MB = total_params * 4 / (1024 ** 2) # float32 => 4 bytes
119
120 print(f"    _Total_Parameters:_{total_params:,}")
121 print(f"    _Trainable_Parameters:_{trainable_params:,}")
122 print(f"    _Estimated_Parameter_Size:_{param_size_MB:.2f}_MB")
123
124 #####
125 # Training Setup #
126 #####
127 history = {"train_loss": [], "val_loss": []}
128 os.makedirs(save_dir, exist_ok=True)
129 save_path = Path(save_dir) / save_name
130
131 best_val_loss = float("inf")
132 best_model_config = None
133 epochs_no_improve = 0
134
135 for epoch in range(num_epochs):
136     # Training Phase
137     model.train()
138     running_loss = 0.0
139
140     for batch in tqdm(train_loader, desc=f"Epoch_{epoch+1}/{num_epochs}_[Train]", unit='_
        Batches'):
141         images = batch['image'].to(device) # (B, 3, H, W)
142         gt_masks = batch['gt_mask'].to(device) # (B, 1, H, W)
143         prompt_heatmaps = batch['prompt_heatmap'].to(device) # (B, 1, H, W)
144         point_classes = batch['point_class'].to(device) # (B,)
145
146         # Create target masks based on point class
147         target_masks = create_point_wise_mask(
148             point_classes,
149             gt_masks
150         )
151
152         optimizer.zero_grad()
153
154         # Forward pass
155         outputs = model(
156             image=images,
157             prompt_heatmap=prompt_heatmaps,
158             point_class=None
159         )
160
161         # Calculate loss
162         loss = train_criterion(outputs, target_masks)
163         loss.backward()
164         optimizer.step()
165
166         running_loss += loss.item()
167
168     train_loss = running_loss / len(train_loader)
169     history["train_loss"].append(train_loss)
170
171     # Validation Phase
172     model.eval()
173     val_loss = 0.0
174
175     with torch.no_grad():
176         for batch in tqdm(val_loader, desc=f"Epoch_{epoch+1}/{num_epochs}_[Val]", unit='_
            Batches'):
177             images = batch['image'].to(device) # (B, 3, H, W)
178             gt_masks = batch['gt_mask'].to(device) # (B, 1, H, W)
179             prompt_heatmaps = batch['prompt_heatmap'].to(device) # (B, 1, H, W)

```

```

180         point_classes = batch['point_class'].to(device) # (B,)
181
182         # Create target masks based on point class
183         target_masks = create_point_wise_mask(
184             point_classes,
185             gt_masks
186         )
187
188         # Forward pass
189         outputs = model(
190             image=images,
191             prompt_heatmap=prompt_heatmaps,
192             point_class=None
193         )
194         loss = val_criterion(outputs, target_masks)
195
196         val_loss += loss.item()
197
198     val_loss /= len(val_loader)
199     history["val_loss"].append(val_loss)
200     print(f"Epoch_{epoch+1}/{num_epochs}_|_Train_Loss:_{train_loss:.6f}_|_Val_Loss:_{val_loss:.6f}")
201
202     # Early Stopping Logic
203     if val_loss < best_val_loss:
204         best_val_loss = val_loss # Update best loss
205         best_epoch = epoch + 1
206         epochs_no_improve = 0 # Reset counter
207         best_model_config = model.state_dict() # record the best model checkpoint so far
208         # Ensure the directory exists
209         tmp_save_dir = Path(save_dir) / "tmp_prompt_checkpoint"
210         tmp_save_dir.mkdir(parents=True, exist_ok=True)
211
212         # Construct save path
213         tmp_save_name = f'best_prompt_checkpoint_epoch_{epoch+1}.pth'
214         tmp_save_path = tmp_save_dir / tmp_save_name
215
216         torch.save(best_model_config, tmp_save_path)
217     else:
218         epochs_no_improve += 1 # Increment counter
219         print(f"No_improvement_for_{epochs_no_improve}_epochs.")
220
221     if epochs_no_improve >= patience:
222         # Save model
223         torch.save(best_model_config, save_path) # Save best model
224         print(f"Early_stopping_triggered_after_{epoch+1}_epochs!")
225         print(f"Best_val_loss_at_epoch_{best_epoch}:_{best_val_loss:.8}")
226         print(f"Best_model_config_saved_at:_{save_path}")
227         return history
228
229     torch.save(best_model_config, save_path) # save the model if the early stop is not triggered
230     print(f"Model_config_saved_at:_{save_path}_ (early_stop_not_triggered)")
231     return history

```

E.2. Training Visualization *training_plot.py*

Generates loss and metric plots across training epochs for both training and validation sets, aiding in model performance analysis and debugging. *training_plot.py*

```

1 import matplotlib.pyplot as plt
2
3 def training_plot(history, save_path=None):
4     """
5     Plots the training and validation loss curves from the history dict.
6
7     Args:

```

```

8     history (dict): Contains 'train_loss' and 'val_loss' lists.
9     save_path (str, optional): If provided, saves the plot to this path.
10    title (str): Title of the plot.
11
12    """
13    plt.figure(figsize=(8, 6))
14    plt.plot(history['train_loss'], label='Train_Loss', marker='o', ms=4)
15    plt.plot(history['val_loss'], label='Validation_Loss', marker='o', ms=4)
16    plt.xlabel('Epoch')
17    plt.legend()
18    plt.grid(True)
19
20    if save_path:
21        plt.savefig(save_path, dpi=300)
22        print(f"Loss_plot_saved_to_{save_path}")
23    else:
24        plt.show()
25
26    plt.close()

```

E.3. Loss Functions *loss_functions.py*

Defines custom loss functions including Focal Loss, Dice Loss, IoU Loss, and combined variants. These are tailored to handle class imbalance and improve segmentation accuracy. *loss_functions.py*

```

1 import torch.nn as nn
2 import torch
3 import torch.nn.functional as F
4 #####
5 ### Loss Functions ###
6 #####
7
8 class DiceLoss(nn.Module):
9     '''
10    calculate the (1 - Dice) loss
11    Args:
12        ouputs(torch.Tensor): model outputs, shape (B, 3, 512, 512)
13        masks(torch.Tensor): ground truth, shape (B, 512, 512)
14    Returns:
15        (1 - mean_dice_over_classes) Loss
16    '''
17    def __init__(self, smooth=1e-6):
18        super(DiceLoss, self).__init__()
19        self.smooth = smooth
20
21    def forward(self, outputs, masks):
22        # Convert logits to probabilities using softmax
23        outputs = F.softmax(outputs, dim=1)
24
25        # Convert masks to one-hot encoding
26        masks_one_hot = F.one_hot(masks.long(), num_classes=outputs.shape[1]).permute(0, 3, 1, 2)
27        .float()
28
29        # Compute Dice Score
30        intersection = torch.sum(outputs * masks_one_hot, dim=(2, 3))
31        union = torch.sum(outputs, dim=(2, 3)) + torch.sum(masks_one_hot, dim=(2, 3))
32
33        dice_score = (2. * intersection + self.smooth) / (union + self.smooth)
34        return 1 - dice_score.mean()
35
36 class IouLoss(nn.Module):
37     '''
38    calculate the (1 - IoU) loss
39    Args:
40        ouputs(torch.Tensor): model outputs, shape (B, 3, 512, 512)
41        masks(torch.Tensor): ground truth, shape (B, 512, 512)
42    Returns:

```

```

42     (1 - mean_iou_over_classes) Loss
43     '''
44     def __init__(self, smooth=1e-6):
45         super().__init__()
46         self.smooth = smooth
47     def forward(self, outputs, masks):
48         # Convert logits to probabilities using softmax
49         outputs = F.softmax(outputs, dim=1)
50
51         # Convert masks to one-hot encoding
52         masks_one_hot = F.one_hot(masks.long(), num_classes=outputs.shape[1]).permute(0, 3, 1, 2)
53         .float()
54
55         # Compute IoU Score
56         intersection = torch.sum(outputs * masks_one_hot, dim=(2, 3))
57         union = torch.sum(outputs + masks_one_hot, dim=(2, 3)) - intersection
58
59         iou_score = (intersection + self.smooth) / (union + self.smooth)
60         return 1 - iou_score.mean()
61
62     def forward(self, outputs, masks):
63         # Convert logits to probabilities using softmax
64         outputs = F.softmax(outputs, dim=1)
65
66         # Convert masks to one-hot encoding
67         masks_one_hot = F.one_hot(masks.long(), num_classes=outputs.shape[1]).permute(0, 3, 1, 2)
68         .float()
69
70         # Compute Dice Score
71         intersection = torch.sum(outputs * masks_one_hot, dim=(2, 3))
72         union = torch.sum(outputs, dim=(2, 3)) + torch.sum(masks_one_hot, dim=(2, 3))
73
74         dice_score = (2. * intersection + self.smooth) / (union + self.smooth)
75         return 1 - dice_score.mean()
76
77 class FocalLoss(nn.Module):
78     """
79     Focal Loss for class imbalance.
80     Args:
81         gamma (float): Focusing parameter (>1 focuses more on hard samples).
82         alpha (tensor, optional): Class weighting (for additional balance).
83         reduction (str): 'mean' or 'sum'.
84     Returns:
85         torch.Tensor: Focal loss.
86     """
87     def __init__(self, gamma=2.0, alpha=None, reduction='mean'):
88         """
89         Args:
90             gamma (float): Focusing parameter (>1 focuses more on hard samples).
91             alpha (tensor, optional): Class weighting (for additional balance).
92             reduction (str): 'mean' or 'sum'.
93         """
94         super().__init__()
95         self.gamma = gamma
96         self.alpha = alpha
97         self.reduction = reduction
98
99     def forward(self, inputs, targets):
100         """
101         Args:
102             inputs (torch.Tensor): Model logits (B, C, H, W).
103             targets (torch.Tensor): Ground truth (B, H, W).
104
105         Returns:
106             torch.Tensor: Focal loss.
107         """
108         log_probs = F.log_softmax(inputs, dim=1) # Convert logits to log-probs
109         probs = torch.exp(log_probs) # Convert to probabilities

```



```

108     # Gather log-probabilities of correct class
109     target_log_probs = log_probs.gather(1, targets.unsqueeze(1)) # (B, 1, H, W)
110     target_probs = probs.gather(1, targets.unsqueeze(1))
111
112     # Compute Focal Loss term
113     focal_weight = (1 - target_probs) ** self.gamma # Focus on hard examples
114
115     if self.alpha is not None:
116         alpha_factor = self.alpha[target_probs] if isinstance(self.alpha, torch.Tensor) else self.alpha
117         focal_weight = focal_weight * alpha_factor
118
119     loss = -focal_weight * target_log_probs # Apply weighting
120
121     if self.reduction == 'mean':
122         return loss.mean()
123     elif self.reduction == 'sum':
124         return loss.sum()
125     else:
126         return loss
127
128 class CombinedFocalDiceLoss(nn.Module):
129     """
130     Combined Loss = (1 - alpha)*FocalLoss + alpha*DiceLoss
131     Args:
132         alpha (float): Balance between Focal and Dice Loss.
133         gamma (float): Focal Loss focusing parameter.
134         focal_alpha (tensor, optional): Class weighting for Focal Loss.
135         smooth (float): Smoothing factor for Dice Loss.
136     Returns:
137         torch.Tensor: Combined loss.
138     """
139     def __init__(self, alpha=0.5, gamma=2.0, focal_alpha=None, smooth=1e-6):
140         super().__init__()
141         self.alpha = alpha
142         self.focal_loss = FocalLoss(gamma=gamma, alpha=focal_alpha)
143         self.dice_loss = DiceLoss(smooth=smooth)
144
145     def forward(self, outputs, masks):
146         fl = self.focal_loss(outputs, masks)
147         dl = self.dice_loss(outputs, masks)
148         return (1 - self.alpha) * fl + self.alpha * dl
149
150 # Binary Loss Function
151 class BinaryFocalLoss(nn.Module):
152     """
153     Focal Loss for binary classification in prompt-based training.
154     Args:
155         alpha (float): Balance between positive and negative classes.
156         gamma (float): Focusing parameter (>1 focuses more on hard samples).
157         reduction (str): 'mean' or 'sum'.
158     Returns:
159         torch.Tensor: Focal loss.
160     """
161     def __init__(self, alpha=0.25, gamma=2.0, reduction='mean'):
162         super().__init__()
163         self.gamma = gamma
164         self.alpha = alpha
165         self.reduction = reduction
166
167     def forward(self, logits, targets):
168         # Apply sigmoid to get probabilities
169         probs = torch.sigmoid(logits)
170         probs = probs.clamp(min=1e-7, max=1 - 1e-7) # avoid log(0)
171
172         # Focal loss formula
173         pt = probs * targets + (1 - probs) * (1 - targets)

```

```

175     alpha_t = self.alpha * targets + (1 - self.alpha) * (1 - targets)
176     focal_loss = -alpha_t * (1 - pt) ** self.gamma * torch.log(pt)
177
178     if self.reduction == 'mean':
179         return focal_loss.mean()
180     elif self.reduction == 'sum':
181         return focal_loss.sum()
182     else:
183         return focal_loss

```

F. Inference & Evaluation

This appendix shows the pipeline of inference and evaluation and related utilities.

F.1. Modular Inference Runner *run_inference.py*

run_inference.py

```

1  #!/usr/bin/env python
2
3  import argparse
4  import torch
5  import os
6  import clip
7  from pathlib import Path
8
9  # import models
10 # Import your models
11 from models.unet_segmentation import UNet
12 from models.autoencoder_segmentation import AutoEncoder, AutoEncoderSegmentation # <-- Import
   your Autoencoder models
13 from models.clip_segmentation import CLIPSegmentationModel
14 from models.prompt_segmentation import PromptSegmentation
15 from utils.inference import inference, promptInference # Ensure you have a separate inference
   function
16
17 # Define function to load the correct model
18 def load_model(mode, checkpoint_path, device, pretrain_path=None):
19     """
20     Load the segmentation model based on the selected mode.
21
22     Args:
23     mode (int): 0 for U-Net, 1 for Autoencoder, 2 for CLIP.
24     checkpoint_path (str): Path to the trained model checkpoint.
25     device (str): Device to use ("cuda" or "cpu").
26
27     Returns:
28     model (torch.nn.Module): The loaded model.
29     """
30
31     #####
32     ##### U-Net #####
33     #####
34     if mode == 0:
35         print("Loading_U-Net_model...")
36         model = UNet()
37         checkpoint = torch.load(checkpoint_path, map_location=device)
38         model.load_state_dict(checkpoint)
39         model.to(device)
40         model.eval()
41         return model
42
43     #####
44     ##### Autoencoder #####
45     #####
46     elif mode == 1:
47         print("Loading_Autoencoder_Segmentation_Model...")

```

```

48
49     if pretrain_path is None:
50         raise ValueError("Pretrained_autoencoder_path_must_be_provided_in_Autoencoder_
51                               mode_(--pretrain_path).")
52
53     # Load pretrained autoencoder
54     autoencoder = AutoEncoder()
55     pretrain_checkpoint = torch.load(pretrain_path, map_location=device)
56     autoencoder.load_state_dict(pretrain_checkpoint)
57     autoencoder.to(device)
58
59     # Wrap the frozen encoder with the segmentation head
60     model = AutoEncoderSegmentation(encoder=autoencoder.encoder, num_classes=3)
61
62     # Load segmentation checkpoint (contains trained decoder)
63     seg_checkpoint = torch.load(checkpoint_path, map_location=device)
64     model.load_state_dict(seg_checkpoint)
65
66     model.to(device)
67     model.eval()
68     return model
69
70 #####
71 ##### CLIP #####
72 #####
73 elif mode == 2:
74     print("Loading_CLIP_Segmentation_model...")
75     clip_model, _ = clip.load("RN50", device=device)
76     model = CLIPSegmentationModel(clip_model=clip_model, num_classes=3)
77     # Load checkpoint
78     checkpoint = torch.load(checkpoint_path, map_location=device)
79     model.load_state_dict(checkpoint)
80     model.to(device)
81     model.float()
82     model.eval()
83     return model
84
85 #####
86 # Prompt-based #
87 #####
88 elif mode == 3:
89     print("Loading_Prompt_Segmentation_model...")
90     model = PromptSegmentation()
91     checkpoint = torch.load(checkpoint_path, map_location=device)
92     model.load_state_dict(checkpoint)
93     model.to(device)
94     model.float()
95     model.eval()
96     return model
97
98 else:
99     raise ValueError("Invalid_mode!_Use_0_for_U-Net,_1_for_Autoencoder,_2_for_CLIP_or_3_for_
100                       prompt_model.")
101
102 # Parse command-line arguments
103 def parse_args():
104     parser = argparse.ArgumentParser(description="Run_inference_for_image_segmentation.")
105     parser.add_argument("--image_dir", type=str, required=True, help="Path_to_input_images")
106     parser.add_argument("--mask_dir", type=str, required=True, help="Path_to_ground_truth_masks")
107     parser.add_argument("--point_dir", type=str, default=None, help="Path_to_prompt_points_(only_
108                               for_prompt_model)")
109     parser.add_argument("--save_dir", type=str, required=True, help="Path_to_save_output_masks")
110     parser.add_argument("--checkpoint_path", type=str, required=True, help="Path_to_model_
111                               checkpoint")
112     parser.add_argument("--pretrain_path", type=str, default=None, help="Path_to_pretrained_
113                               autoencoder_(required_for_autoencoder_mode)")
114     parser.add_argument("--target_size", type=int, default=512, help="Target_input_size_for_model
115                               ")

```

```

110     parser.add_argument("--device", type=str, required=True, help="Device_for_inference_(cuda/cpu
111         )")
112     parser.add_argument("--mode", type=int, required=True, help="Model_selection:_0_for_U-Net,_1_
113         for_Autoencoder,_2_for_CLIP,_3_for_prompt_model")
114
115     return parser.parse_args()
116
117 if __name__ == "__main__":
118     args = parse_args()
119
120     # Ensure save directory exists
121     os.makedirs(args.save_dir, exist_ok=True)
122
123     # Load model
124     model = load_model(args.mode, args.checkpoint_path, args.device, pretrain_path=args.
125         pretrain_path)
126     # Run inference
127     if args.mode in [0,1,2]:
128         inference(
129             image_dir=args.image_dir,
130             mask_dir=args.mask_dir,
131             save_dir=args.save_dir,
132             mode = args.mode,
133             input_image_size=args.target_size,
134             model=model,
135             device=args.device
136         )
137     elif args.mode == 3:
138         promptInference(
139             image_dir=args.image_dir,
140             mask_dir=args.mask_dir,
141             point_dir= args.point_dir,
142             gt_save_dir= Path(args.save_dir) / "gt/",
143             pred_save_dir=Path(args.save_dir) / "pred/",
144             threshold=0.5,
145             input_image_size=args.target_size,
146             model=model,
147             device=args.device
148         )
149     else:
150         raise ValueError("Invalid_mode!_Use_0_for_U-Net,_1_for_Autoencoder,_2_for_CLIP_or_3_for_
151             prompt_model.")

```

F.2. Inference Blocks *inference.py*

inference.py

```

1 # run_inference.py
2
3 import torch
4 from data.preprocessing import (class2color,
5     resize_with_padding,
6     clip_transform,
7     standard_transform,
8     create_point_wise_mask)
9
10 import os
11 from utils.restore_mask_size import restore_original_mask
12 from data.PetDataset import PetDataset, PetDatasetWithPrompt
13 from torch.utils.data import DataLoader
14 from pathlib import Path
15 from tqdm import tqdm
16 import numpy as np
17 from PIL import Image
18
19 def custom_collate_fn(batch):
20     images, masks, initial_img_sizes, filenames = zip(*batch) # Unpack batch
21
22     # Convert images and masks to tensors (PyTorch default behavior)

```

```

22 images = torch.stack(images, dim=0)
23 masks = torch.stack(masks, dim=0)
24
25 return images, masks, list(initial_img_sizes), list(filenamees) # Keep 'initial_img_sizes' as
    list of tuples
26
27 def inference(image_dir: str, mask_dir: str, save_dir: str, input_image_size: int,
28               mode: int, model, device: str) -> np.ndarray:
29     """
30     Perform inference on input images and save the predicted segmentation masks.
31
32     Args:
33         image_dir (str): Directory containing input images.
34         mask_dir (str): Directory containing ground truth masks (for size reference).
35         save_dir (str): Directory to save the output predicted masks.
36         input_image_size (int): Target size for resizing (used in CLIP mode).
37         mode (int): Model mode selector:
38             0 - U-Net,
39             1 - Autoencoder,
40             2 - CLIP.
41         model (torch.nn.Module): Pre-loaded segmentation model ready for inference.
42         device (str): Computation device ("cuda" or "cpu").
43
44     Returns:
45         None. (Predicted color masks are saved to 'save_dir')
46     """
47
48     # Initialize paths of images (list of str)
49     image_paths = sorted(Path(image_dir).glob("*.*)" )
50     mask_paths = sorted(Path(mask_dir).glob("*.*)" )
51
52     # assign the correct torch dataset format
53     if mode == 0 or mode == 1:
54         test_dataset = PetDataset(
55             img_paths = image_paths,
56             msk_paths = mask_paths,
57             resize_fn= resize_with_padding,
58             resize_target_size= 512,
59             transform = standard_transform)
60     elif mode == 2:
61         test_dataset = PetDataset(
62             img_paths = image_paths,
63             msk_paths = mask_paths,
64             resize_fn = resize_with_padding,
65             resize_target_size = input_image_size,
66             transform = clip_transform
67         )
68
69     else:
70         raise ValueError("Invalid mode. Use 0 for U-Net, 1 for autoencoder-based segmentation or 2 for CLIP-based segmentation.")
71
72     test_loader = DataLoader(test_dataset, batch_size=4, num_workers=4, collate_fn=
        custom_collate_fn)
73
74
75     with torch.no_grad():
76         total_batches = len(test_loader) # Get total number of batches
77         for images, _, initial_img_sizes, filenamees in \
78             tqdm(test_loader, desc=f"Processing_{len(test_loader.dataset)}_images_in_{total_batches}_
                batches"):
79             images = images.to(device)
80             outputs = model(images)
81
82             # Convert logits to class labels (B, C, H, W) (B, H, W)
83             pred_masks = torch.argmax(outputs, dim=1).cpu().numpy()
84             resized_pred_masks = [restore_original_mask(pred_masks[i], initial_img_sizes[i]) for
                i in range(len(images))]

```

```

85
86     # Save predicted masks
87     for i, filename in enumerate(filenamees):
88         pred_mask_img = class2color(resized_pred_masks[i]) # Convert class labels to
            color mask
89         # Ensure data type is uint8 to avoid artifacts
90         pred_mask_img = np.array(pred_mask_img, dtype=np.uint8)
91         save_img = Image.fromarray(pred_mask_img, mode="RGB") # Convert the ndarray into
            RGB Image
92         save_path = os.path.join(save_dir, Path(filename).stem+'.png') # Keep original
            filename
93         save_img.save(save_path, format='PNG')
94
95     # print the dst directory
96     print(f"Predicted_masks_saved_to_{save_dir}")
97
98 def promptInference(image_dir: str, mask_dir:str, point_dir:str,
99                     gt_save_dir:str, pred_save_dir:str,
100                     threshold:float, input_image_size: int,
101                     model: callable, device: str) -> np.ndarray:
102     '''
103     Run the inference on test images with input of prompt points
104
105     Args:
106     image_dir (str): Directory containing input images
107     mask_dir (str): Directory containing ground truth mask images
108     point_dir (str): Directory containing point prompts
109     gt_save_dir (str): Directory where ground truth masks will be saved
110     pred_save_dir (str): Directory where predicted masks will be saved
111     threshold (float): Threshold value for binary segmentation (sigmoid output > threshold)
112     input_image_size (int): Target size for resizing images before inference
113     model (callable): The segmentation model to use for inference
114     device (str): Device to run inference on (e.g., 'cuda', 'cpu')
115
116     Returns:
117     np.ndarray: Processed data array containing segmentation results
118
119     Note:
120     - The function processes batches of images with corresponding masks and point prompts
121     - Predicted masks are saved as PNG files in pred_save_dir
122     - Ground truth masks are saved as PNG files in gt_save_dir
123     - File names are preserved from the original image files
124     '''
125     image_paths = sorted(Path(image_dir).glob("*..*"))
126     mask_paths = sorted(Path(mask_dir).glob("*..*"))
127     point_paths = sorted(Path(point_dir).glob("*..*"))
128
129     test_dataset = PetDatasetWithPrompt(
130         img_paths=image_paths,
131         msk_paths=mask_paths,
132         pnt_paths=point_paths,
133         resize_fn=resize_with_padding,
134         resize_target_size=input_image_size,
135         transform=standard_transform,
136         load_multiple_points=True
137     )
138
139     test_loader = DataLoader(test_dataset, batch_size=4, num_workers=4, collate_fn=
        prompt_custom_collate_fn)
140
141
142     with torch.no_grad():
143         total_batches = len(test_loader) # Get total number of batches
144         for batch in tqdm(test_loader, desc=f"Processing_{len(test_loader.dataset)}_images_in_{
            total_batches}_batches"):
145             images = batch['image'].to(device) # (B, 3, H, W)
146             gt_masks = batch['gt_mask'].to(device) # (B, 1, H, W)
147             prompt_heatmaps = batch['prompt_heatmap'].to(device) # (B, 1, H, W)

```

```

148     point_classes = batch['point_class'].to(device)          # (B,)
149     initial_img_sizes = batch['initial_img_size']
150     filenames = batch['img_path']
151     # Create target masks based on point class
152     target_masks = create_point_wise_mask(
153         point_classes,
154         gt_masks
155     )
156
157     with torch.no_grad():
158         output = model(image=images,
159                         prompt_heatmap=prompt_heatmaps,
160                         point_class = point_classes)
161
162     # Convert logits to class labels (B, C, H, W)          (B, H, W)
163     pred_masks = torch.sigmoid(output) > threshold
164     # restore the gt and pred masks shape
165     pred_masks = pred_masks.squeeze(1).cpu().numpy().astype(np.uint8)
166     target_masks = target_masks.squeeze(1).cpu().numpy().astype(np.uint8)
167
168     resized_pred_masks = [restore_original_mask(pred_masks[i], initial_img_sizes[i][:2])
169                           for i in range(len(images))]
170     resized_gt_masks = [restore_original_mask(target_masks[i], initial_img_sizes[i][:2])
171                         for i in range(len(images))]
172     # Save prompt gt binary mask
173     for i, filename in enumerate(filenames):
174         mask_arr = (resized_gt_masks[i]*255).astype(np.uint8)
175         mask_save = Image.fromarray(mask_arr, mode='L')
176         mask_save_path = os.path.join(gt_save_dir, Path(filename).stem+'.png') # Keep
177         original filename
178         # Ensure parent directory exists
179         Path(mask_save_path).parent.mkdir(parents=True, exist_ok=True)
180         mask_save.save(mask_save_path, format='PNG')
181
182     # Save predicted masks
183     for i, filename in enumerate(filenames):
184         pred_arr = (resized_pred_masks[i]*255).astype(np.uint8)
185         pred_save = Image.fromarray(pred_arr, mode="L") # Convert the ndarray into RGB
186         Image
187         pred_save_path = os.path.join(pred_save_dir, Path(filename).stem+'.png') # Keep
188         original filename
189         # Ensure parent directory exists
190         Path(pred_save_path).parent.mkdir(parents=True, exist_ok=True)
191         pred_save.save(pred_save_path, format='PNG')
192
193 def prompt_custom_collate_fn(batch):
194     # batch is a list of dictionaries, one per sample
195     # We'll build batched outputs manually
196
197     images = torch.stack([sample['image'] for sample in batch], dim=0)
198     gt_masks = torch.stack([sample['gt_mask'] for sample in batch], dim=0)
199     prompt_heatmaps = torch.stack([sample['prompt_heatmap'] for sample in batch], dim=0)
200     point_classes = torch.stack([sample['point_class'] for sample in batch], dim=0)
201     # Keep 'initial_img_size' as a list instead of stacking into a Tensor
202     initial_img_sizes = [sample['initial_img_size'] for sample in batch]
203     img_paths = [sample['img_path'] for sample in batch]
204
205     return {
206         'image': images,
207         'gt_mask': gt_masks,
208         'prompt_heatmap': prompt_heatmaps,
209         'point_class': point_classes,
210         'initial_img_size': initial_img_sizes,
211         'img_path': img_paths
212     }

```

F.3. Mask Size Restoration *restore_mask_size.py*

It restores the predicted masks to the initial sizes for the fair calculation of IoU *restore_mask_size.py*

```
1 import cv2 as cv
2 import numpy as np
3
4 def restore_original_mask(mask: np.ndarray, original_size: tuple) -> np.ndarray:
5     """
6     Resize the class mask back to the original image size.
7
8     Args:
9         mask (np.ndarray): The resized mask of shape (target_size, target_size).
10        original_size (tuple): The original image dimensions (H, W).
11
12    Returns:
13        np.ndarray: The restored mask with original dimensions (H, W).
14    """
15    target_size = mask.shape[0] # Assuming square mask (224, 224)
16    orig_h, orig_w = original_size
17
18    # Compute scale factor (same used in resize_with_padding)
19    scale = target_size / max(orig_w, orig_h)
20    new_w, new_h = int(orig_w * scale), int(orig_h * scale)
21
22    # Compute padding offsets
23    paste_x = (target_size - new_w) // 2
24    paste_y = (target_size - new_h) // 2
25
26    # Crop the valid region (remove padding)
27    cropped_mask = mask[paste_y:paste_y + new_h, paste_x:paste_x + new_w]
28
29    # Resize back to original dimensions using NEAREST to keep discrete labels
30    original_mask = cv.resize(cropped_mask, (orig_w, orig_h), interpolation=cv.INTER_NEAREST)
31
32    return original_mask
```

F.4. IoU Computation *compute_IoU.py*

computer_IoU.py

```
1 import argparse
2 from pathlib import Path
3 from PIL import Image
4 from tqdm import tqdm
5 import numpy as np
6 from data.preprocessing import color2class
7 import torch
8
9 def binarize_mask(mask_rgb):
10     """
11     Converts an RGB mask to binary (1: foreground, 0: background)
12     Args:
13         mask_rgb (np.ndarray): The RGB mask image.
14     Returns:
15         np.ndarray: Binary mask where 1 represents the foreground and 0 the background.
16     """
17     gray = np.array(Image.fromarray(mask_rgb).convert("L"))
18     binary = (gray > 127).astype(np.uint8)
19     return binary
20
21 def compute_dataset_iou(mode:str, gt_folder:str, pred_folder:str, output_file:str) -> None:
22     """
23     Computes the Intersection over Union (IoU) for a dataset of segmentation masks.
24     Args:
25         mode (str): Mode of evaluation, either "3-class" or "binary".
26         gt_folder (str): Path to the ground truth masks folder.
27         pred_folder (str): Path to the predicted masks folder.
```



```

28     output_file (str): Path to save the IoU results.
29 Returns:
30     None. (results saved to <output_file>.txt)
31 '''
32 print(f'Calculating IoU of predicted images in {pred_folder} at mode {mode}...')
33
34 # Define dataset paths
35 gt_paths = sorted(Path(gt_folder).glob("*.png"))
36 pred_paths = sorted(Path(pred_folder).glob("*.png"))
37 device = 'cuda' if torch.cuda.is_available() else 'cpu'
38
39 # Initialize accumulators
40 if mode == "3-class":
41     intersection_sum = {0: 0, 1: 0, 2: 0}
42     union_sum = {0: 0, 1: 0, 2: 0}
43 else: # binary
44     intersection_sum = {1: 0} # foreground only
45     union_sum = {1: 0}
46
47 # Open file in write mode
48 with open(output_file, "w") as file:
49     for idx in tqdm(range(len(pred_paths)), desc='Test_Image_Loading', unit='image'):
50         if mode == "3-class":
51             gt_mask = Image.open(gt_paths[idx]).convert('RGB')
52             gt_mask = np.array(gt_mask)
53             gt_classes = color2class(gt_mask)
54             gt_classes = torch.from_numpy(gt_classes).to(device).long()
55
56             pred_mask = Image.open(pred_paths[idx]).convert('RGB')
57             pred_mask = np.array(pred_mask)
58             pred_classes = color2class(pred_mask)
59             pred_classes = torch.from_numpy(pred_classes).to(device).long()
60
61             # convert pred and gt images into one-hot embedding
62             pred_one_hot = torch.nn.functional.one_hot(pred_classes, num_classes=3).permute(
63                 (2,0,1)).float()
64             gt_one_hot = torch.nn.functional.one_hot(gt_classes, num_classes=3).permute(2, 0,
65                 1).float()
66             intersection = (pred_one_hot * gt_one_hot).sum(dim=(1, 2))
67             union = pred_one_hot.sum(dim=(1, 2)) + gt_one_hot.sum(dim=(1, 2)) - intersection
68
69             for cls in range(3):
70                 intersection_sum[cls] += intersection[cls].item()
71                 union_sum[cls] += union[cls].item()
72         # For binary mode, we only care about the foreground (1) and background (0)
73         elif mode == "binary":
74             gt_mask = Image.open(gt_paths[idx]).convert('RGB')
75             pred_mask = Image.open(pred_paths[idx]).convert('RGB')
76
77             gt_binary = binarize_mask(np.array(gt_mask))
78             pred_binary = binarize_mask(np.array(pred_mask))
79
80             gt_tensor = torch.from_numpy(gt_binary).to(device)
81             pred_tensor = torch.from_numpy(pred_binary).to(device)
82
83             intersection = torch.logical_and(pred_tensor, gt_tensor).sum().item()
84             union = torch.logical_or(pred_tensor, gt_tensor).sum().item()
85
86             intersection_sum[1] += intersection
87             union_sum[1] += union
88         # initialize final iou dict
89         final_iou = {}
90         # Compute IoU per class
91         if mode == "3-class":
92             for cls in range(3):
93                 if union_sum[cls] > 0:
94                     final_iou[cls] = intersection_sum[cls] / union_sum[cls]
95                 else:

```

```

94         final_iou[cls] = None # Optional: mark as N/A if class never appears
95
96     # Mean IoU: Only include classes that exist
97     valid_iou = [iou for iou in final_iou.values() if iou is not None]
98     mean_iou = sum(valid_iou) / len(valid_iou) if valid_iou else 0
99
100    # Final results
101    final_results = (f"\nFinal_Dataset-level_IoU_Results:\n"
102                    f"IoU_of_Background:_{final_iou[0]:.4f}\n"
103                    f"IoU_of_Cats:_{final_iou[1]:.4f}\n"
104                    f"IoU_of_Dogs:_{final_iou[2]:.4f}\n"
105                    f"Mean_IoU:_{mean_iou:.4f}\n")
106
107    elif mode == "binary":
108        if union_sum[1] > 0:
109            iou = intersection_sum[1] / union_sum[1]
110        else:
111            iou = 0.0
112
113        final_results = (f"\nFinal_Dataset-level_IoU_Results_(Binary):\n"
114                        f"IoU_of_Foreground_vs_Background:_{iou:.4f}\n")
115    else:
116        raise ValueError("Invalid_mode._Choose_'3-class'_or_'binary'.")
117    # Save results to file
118    print(final_results.strip())
119    file.write(final_results)
120
121    print(f"\nDataset-level_IoU_results_saved_in:_{output_file}")
122
123 if __name__ == "__main__":
124     parser = argparse.ArgumentParser(description="Compute_Dataset-level_IoU_for_Segmentation_Masks")
125     parser.add_argument('--mode', type=str, choices=["3-class", "binary"], default="3-class",
126                         help='Evaluation_mode:_3-class_(default)_or_binary_(SAM-style_fg/bg)')
127     parser.add_argument('--gt_folder', type=str, required=True, help='Path_to_the_ground_truth_masks_folder')
128     parser.add_argument('--pred_folder', type=str, required=True, help='Path_to_the_predicted_masks_folder')
129     parser.add_argument('--output_file', type=str, default='iou_results.txt', help='Path_to_save_the_IoU_results')
130
131     args = parser.parse_args()
132     compute_dataset_iou(args.mode, args.gt_folder, args.pred_folder, args.output_file)

```

G. UI

This appendix presents the implementation of the user interface (UI) used for prompt-based segmentation. The UI allows users to interactively provide input prompts—such as points or boxes—used to guide the segmentation model.

G.1. Source Code

```

1 import sys
2 import os
3 import numpy as np
4 import torch
5 from PIL import Image
6 import cv2
7 from pathlib import Path
8 from PyQt5.QtWidgets import (QApplication, QMainWindow, QWidget, QVBoxLayout, QHBoxLayout,
9                             QPushButton, QFileDialog, QLabel, QRadioButton, QButtonGroup,
10                             QSlider, QGroupBox, QComboBox)
11 from PyQt5.QtGui import QPixmap, QPainter, QColor, QPen, QImage
12 from PyQt5.QtCore import Qt, QPoint, QRect
13
14 from models.prompt_segmentation import PromptSegmentation
15 from data.preprocessing import resize_with_padding

```

```

16 from utils.restore_image_size import restore_original_mask
17
18 class SegmentationUI(QMainWindow):
19     def __init__(self):
20         super().__init__()
21         self.setWindowTitle("Interactive_Segmentation_Tool")
22         self.setGeometry(100, 100, 1200, 800)
23
24         # Initialize state variables
25         self.original_image = None
26         self.display_image = None
27         self.original_size = None
28         self.current_mask = None
29         self.combined_mask = None
30         self.prompts = [] # Store points/boxes for inference
31         self.temp_box = None # For box drawing
32         self.mode = "point" # Default mode is point placement
33         self.current_class = 1 # Default class is foreground (1)
34         self.target_size = 512 # Model input size
35         self.threshold = 0.5 # Default threshold for segmentation
36
37         # Load model
38         self.model = self.load_model()
39
40         # Set up UI
41         self.setup_ui()
42
43     def load_model(self):
44         """Load the prompt-based segmentation model"""
45         device = "cuda" if torch.cuda.is_available() else "cpu"
46         print(f"Using_device:_{device}")
47
48         # Load model
49         model = PromptSegmentation()
50         checkpoint_path = r"E:\Oxford-IIIT_Pet_Dataset_Segmentation_test\params\
51             best_prompt_checkpoint_epoch_1.pth"
52         checkpoint = torch.load(checkpoint_path, map_location=device)
53         model.load_state_dict(checkpoint)
54         model.to(device)
55         model.eval()
56
57         return model
58
59     def setup_ui(self):
60         """Set up the user interface"""
61         # Main widget and layout
62         main_widget = QWidget()
63         main_layout = QHBoxLayout()
64
65         # Left side: image display
66         self.canvas = QLabel()
67         self.canvas.setMinimumSize(512, 512)
68         self.canvas.setAlignment(Qt.AlignCenter)
69         self.canvas.setStyleSheet("background-color:_{#f0f0f0};_border:_{1px_solid_#ccc;}")
70         self.canvas.mousePressEvent = self.canvas_click
71         self.canvas.mouseMoveEvent = self.canvas_move
72         self.canvas.mouseReleaseEvent = self.canvas_release
73
74         # Right side: controls
75         controls_layout = QVBoxLayout()
76
77         # Image loading group
78         image_group = QGroupBox("Image")
79         image_layout = QVBoxLayout()
80         self.load_image_btn = QPushButton("Load_Image")
81         self.load_image_btn.clicked.connect(self.load_image)
82         image_layout.addWidget(self.load_image_btn)
83         image_group.setLayout(image_layout)

```

```

83     controls_layout.addWidget(image_group)
84
85     # Prompt tools group
86     prompt_group = QGroupBox("Prompt_Tools")
87     prompt_layout = QVBoxLayout()
88
89     # Prompt type selection
90     prompt_type_layout = QHBoxLayout()
91     self.point_radio = QRadioButton("Point")
92     self.box_radio = QRadioButton("Box")
93     self.point_radio.setChecked(True)
94     self.point_radio.toggled.connect(lambda: self.set_mode("point"))
95     self.box_radio.toggled.connect(lambda: self.set_mode("box"))
96     prompt_type_layout.addWidget(self.point_radio)
97     prompt_type_layout.addWidget(self.box_radio)
98     prompt_layout.addLayout(prompt_type_layout)
99
100    # # Class selection
101    # class_layout = QHBoxLayout()
102    # class_layout.addWidget(QLabel("Class:"))
103    # self.class_combo = QComboBox()
104    # self.class_combo.addItem("Background (0)", "Foreground (1)")
105    # self.class_combo.setCurrentIndex(1) # Default to foreground
106    # self.class_combo.currentIndexChanged.connect(self.change_class)
107    # class_layout.addWidget(self.class_combo)
108    # prompt_layout.addLayout(class_layout)
109
110    # Threshold slider
111    threshold_layout = QHBoxLayout()
112    threshold_layout.addWidget(QLabel("Threshold:"))
113    self.threshold_slider = QSlider(Qt.Horizontal)
114    self.threshold_slider.setMinimum(1)
115    self.threshold_slider.setMaximum(99)
116    self.threshold_slider.setValue(int(self.threshold * 100))
117    self.threshold_slider.valueChanged.connect(self.change_threshold)
118    threshold_layout.addWidget(self.threshold_slider)
119    self.threshold_label = QLabel(f"{self.threshold:.2f}")
120    threshold_layout.addWidget(self.threshold_label)
121    prompt_layout.addLayout(threshold_layout)
122
123    # Action buttons
124    self.run_btn = QPushButton("Run_Segmentation")
125    self.run_btn.clicked.connect(self.run_segmentation)
126    self.run_btn.setEnabled(False)
127    prompt_layout.addWidget(self.run_btn)
128
129    self.clear_prompts_btn = QPushButton("Clear_Prompts")
130    self.clear_prompts_btn.clicked.connect(self.clear_prompts)
131    self.clear_prompts_btn.setEnabled(False)
132    prompt_layout.addWidget(self.clear_prompts_btn)
133
134    self.clear_all_btn = QPushButton("Clear_All")
135    self.clear_all_btn.clicked.connect(self.clear_all)
136    self.clear_all_btn.setEnabled(False)
137    prompt_layout.addWidget(self.clear_all_btn)
138
139    prompt_group.setLayout(prompt_layout)
140    controls_layout.addWidget(prompt_group)
141
142    # Result group
143    result_group = QGroupBox("Results")
144    result_layout = QVBoxLayout()
145
146    self.save_mask_btn = QPushButton("Save_Mask")
147    self.save_mask_btn.clicked.connect(self.save_mask)
148    self.save_mask_btn.setEnabled(False)
149    result_layout.addWidget(self.save_mask_btn)
150

```

```

151     self.toggle_view_btn = QPushButton("Toggle_Overlay")
152     self.toggle_view_btn.clicked.connect(self.toggle_overlay)
153     self.toggle_view_btn.setEnabled(False)
154     result_layout.addWidget(self.toggle_view_btn)
155
156     result_group.setLayout(result_layout)
157     controls_layout.addWidget(result_group)
158
159     # Status indicator
160     self.status_label = QLabel("Ready")
161     controls_layout.addWidget(self.status_label)
162
163     # Add stretcher to push everything up
164     controls_layout.addStretch()
165
166     # Arrange layouts
167     main_layout.addWidget(self.canvas, 3)
168     main_layout.addLayout(controls_layout, 1)
169
170     main_widget.setLayout(main_layout)
171     self.setCentralWidget(main_widget)
172
173     def set_mode(self, mode):
174         """Set the prompt mode (point or box)"""
175         self.mode = mode
176         self.temp_box = None
177         self.update_canvas()
178
179     def change_class(self, index):
180         """Update the current class when combo box selection changes"""
181         self.current_class = 0 if index == 0 else 1
182
183     def change_threshold(self, value):
184         """Update threshold value from slider"""
185         self.threshold = value / 100.0
186         self.threshold_label.setText(f"{self.threshold:.2f}")
187         # Re-run segmentation if we have a current mask
188         if self.current_mask is not None:
189             self.run_segmentation()
190
191     def load_image(self):
192         """Load an image from file"""
193         file_path, _ = QFileDialog.getOpenFileName(
194             self, "Open_Image", "", "Image_Files (*.png *.jpg *.jpeg *.bmp)"
195         )
196
197         if file_path:
198             # Load and store the original image
199             self.original_image = np.array(Image.open(file_path).convert("RGB"))
200             self.original_size = self.original_image.shape[:2]
201
202             # Resize for display
203             self.resize_and_display_image()
204
205             # Enable UI elements
206             self.run_btn.setEnabled(True)
207             self.clear_all_btn.setEnabled(True)
208
209             # Clear existing prompts and masks
210             self.prompts = []
211             self.current_mask = None
212             self.combined_mask = None
213             self.status_label.setText("Image_loaded._Add_prompts_and_run_segmentation.")
214
215     def resize_and_display_image(self):
216         """Resize the original image and update display"""
217         # Resize while maintaining aspect ratio
218         img_resized = resize_with_padding(self.original_image, self.target_size, False)

```

```

219     self.display_image = img_resized.copy()
220
221     # Update canvas
222     self.update_canvas()
223
224     def update_canvas(self):
225         """Update the canvas with current image and overlays"""
226         if self.display_image is None:
227             return
228
229         # Create a copy of the display image for drawing
230         canvas_img = self.display_image.copy()
231
232         # Draw mask overlay if available
233         if self.combined_mask is not None:
234             # Create a colored overlay (semi-transparent)
235             overlay = np.zeros_like(canvas_img)
236             overlay[self.combined_mask > 0] = [0, 255, 0] # Green overlay
237
238             # Blend with original image
239             alpha = 0.5
240             canvas_img = cv2.addWeighted(overlay, alpha, canvas_img, 1-alpha, 0)
241
242         # Draw prompts
243         for prompt in self.prompts:
244             if prompt[0] == "point":
245                 x, y, cls = prompt[1]
246                 color = (0, 0, 255) if cls == 0 else (255, 0, 0) # Red for foreground, blue for
                    background
247                 cv2.circle(canvas_img, (x, y), 5, color, -1)
248             elif prompt[0] == "box":
249                 x1, y1, x2, y2, cls = prompt[1]
250                 color = (0, 0, 255) if cls == 0 else (255, 0, 0)
251                 cv2.rectangle(canvas_img, (x1, y1), (x2, y2), color, 2)
252
253         # Draw temporary box if in progress
254         if self.temp_box is not None:
255             x1, y1, x2, y2 = self.temp_box
256             cv2.rectangle(canvas_img, (x1, y1), (x2, y2), (255, 255, 0), 2)
257
258         # Convert to QImage and display
259         h, w, c = canvas_img.shape
260         q_img = QImage(canvas_img.data, w, h, w*c, QImage.Format_RGB888)
261         self.canvas.setPixmap(QPixmap.fromImage(q_img))
262
263     def canvas_click(self, event):
264         """Handle mouse click on canvas"""
265         if self.display_image is None:
266             return
267
268         # Calculate image position within the canvas
269         canvas_width = self.canvas.width()
270         canvas_height = self.canvas.height()
271         img_height, img_width = self.display_image.shape[:2]
272
273         # Calculate offsets for centered image
274         x_offset = max(0, (canvas_width - img_width) // 2)
275         y_offset = max(0, (canvas_height - img_height) // 2)
276
277         # Adjust coordinates
278         x = event.x() - x_offset
279         y = event.y() - y_offset
280
281         # Check if click is within image bounds
282         if 0 <= x < img_width and 0 <= y < img_height:
283             if self.mode == "point":
284                 # Add point prompt immediately
285                 self.prompts.append(("point", (x, y, self.current_class)))

```

```

286         self.clear_prompts_btn.setEnabled(True)
287         self.update_canvas()
288         elif self.mode == "box":
289             # Start box drawing
290             self.temp_box = [x, y, x, y]
291     else:
292         # Click outside the image area
293         self.status_label.setText("Click_within_the_image_area")
294
295     def canvas_move(self, event):
296         """Handle mouse movement on canvas (for box drawing)"""
297         if self.temp_box is None:
298             return
299
300         # Calculate image position within the canvas
301         canvas_width = self.canvas.width()
302         canvas_height = self.canvas.height()
303         img_height, img_width = self.display_image.shape[:2]
304
305         # Calculate offsets for centered image
306         x_offset = max(0, (canvas_width - img_width) // 2)
307         y_offset = max(0, (canvas_height - img_height) // 2)
308
309         # Adjust coordinates
310         x = event.x() - x_offset
311         y = event.y() - y_offset
312
313         # Constrain coordinates to image bounds
314         x = max(0, min(x, img_width - 1))
315         y = max(0, min(y, img_height - 1))
316
317         self.temp_box[2], self.temp_box[3] = x, y
318         self.update_canvas()
319
320     def canvas_release(self, event):
321         """Handle mouse release on canvas (for box drawing)"""
322         if self.temp_box is None:
323             return
324
325         # Calculate image position within the canvas
326         canvas_width = self.canvas.width()
327         canvas_height = self.canvas.height()
328         img_height, img_width = self.display_image.shape[:2]
329
330         # Calculate offsets for centered image
331         x_offset = max(0, (canvas_width - img_width) // 2)
332         y_offset = max(0, (canvas_height - img_height) // 2)
333
334         # Adjust coordinates
335         x = event.x() - x_offset
336         y = event.y() - y_offset
337
338         # Constrain coordinates to image bounds
339         x = max(0, min(x, img_width - 1))
340         y = max(0, min(y, img_height - 1))
341
342         x1, y1 = self.temp_box[0], self.temp_box[1]
343         x2, y2 = x, y
344
345         # Ensure correct ordering (x1,y1 is top-left, x2,y2 is bottom-right)
346         if x1 > x2:
347             x1, x2 = x2, x1
348         if y1 > y2:
349             y1, y2 = y2, y1
350
351         # Only add if box has some area
352         if x2 > x1 and y2 > y1:
353             self.prompts.append(("box", (x1, y1, x2, y2, self.current_class)))

```

```

354         self.clear_prompts_btn.setEnabled(True)
355
356         self.temp_box = None
357         self.update_canvas()
358
359     def generate_heatmap(self, prompt, img_shape):
360         """Generate a heatmap for the given prompt"""
361         h, w = img_shape
362         heatmap = np.zeros((h, w), dtype=np.float32)
363
364         if prompt[0] == "point":
365             x, y, _ = prompt[1]
366             # Create a Gaussian heatmap centered at the point
367             sigma = 10.0 # Controls the spread of the Gaussian
368             y_grid, x_grid = np.mgrid[0:h, 0:w]
369             heatmap = np.exp(-((x_grid - x) ** 2 + (y_grid - y) ** 2) / (2 * sigma ** 2))
370         elif prompt[0] == "box":
371             x1, y1, x2, y2, _ = prompt[1]
372             # Create a binary heatmap for the box region
373             heatmap[y1:y2+1, x1:x2+1] = 1.0
374
375         return heatmap
376
377     def run_segmentation(self):
378         """Run segmentation based on current prompts"""
379         if not self.prompts or self.display_image is None:
380             self.status_label.setText("Add_at_least_one_prompt_before_running_segmentation")
381             return
382
383         self.status_label.setText("Running_segmentation...")
384
385         # Process each prompt
386         self.combined_mask = None
387
388         for prompt in self.prompts:
389             # Generate heatmap for this prompt
390             heatmap = self.generate_heatmap(prompt, (self.target_size, self.target_size))
391
392             # Prepare inputs for the model - APPLY STANDARD TRANSFORM
393             from data.preprocessing import standard_transform
394             # Convert to PIL Image for transform
395             pil_img = Image.fromarray(self.display_image)
396             # Apply the same transform used during training/inference
397             img_tensor = standard_transform(pil_img).unsqueeze(0) # Add batch dimension
398
399             heatmap_tensor = torch.from_numpy(heatmap).unsqueeze(0).unsqueeze(0).float()
400
401             # Move tensors to device
402             device = next(self.model.parameters()).device
403             img_tensor = img_tensor.to(device)
404             heatmap_tensor = heatmap_tensor.to(device)
405
406             # Run inference
407             with torch.no_grad():
408                 output = self.model(
409                     image=img_tensor,
410                     prompt_heatmap=heatmap_tensor,
411                     point_class= None
412                 )
413
414             # Apply threshold to get binary mask
415             pred_mask = (torch.sigmoid(output) > self.threshold).squeeze().cpu().numpy().astype(
416                 np.uint8)
417
418             # Union with previous masks
419             if self.combined_mask is None:
420                 self.combined_mask = pred_mask
421             else:

```



```

421         self.combined_mask = np.logical_or(self.combined_mask, pred_mask).astype(np.uint8
422         )
423
424         # Enable save mask button
425         self.save_mask_btn.setEnabled(True)
426         self.toggle_view_btn.setEnabled(True)
427
428         # Update canvas with mask overlay
429         self.update_canvas()
430
431         self.status_label.setText("Segmentation_complete")
432
433     def toggle_overlay(self):
434         """Toggle between showing the mask overlay and original image"""
435         if self.combined_mask is None:
436             return
437
438         # Toggle by temporarily removing and restoring the mask
439         if hasattr(self, '_temp_mask'):
440             self.combined_mask = self._temp_mask
441             delattr(self, '_temp_mask')
442         else:
443             self._temp_mask = self.combined_mask
444             self.combined_mask = None
445
446         self.update_canvas()
447
448     def save_mask(self):
449         """Save the current mask to a file"""
450         if self.combined_mask is None:
451             return
452
453         file_path, _ = QFileDialog.getSaveFileName(
454             self, "Save_Mask", "", "PNG_Files_(*.png)"
455         )
456
457         if file_path:
458             # Resize mask back to original image size
459             mask_resized = cv2.resize(
460                 self.combined_mask * 255, # Scale to 0-255
461                 (self.original_size[1], self.original_size[0]),
462                 interpolation=cv2.INTER_NEAREST
463             )
464
465             # Save mask
466             cv2.imwrite(file_path, mask_resized)
467             self.status_label.setText(f"Mask_saved_to_{file_path}")
468
469     def clear_prompts(self):
470         """Clear all prompts but keep the loaded image"""
471         self.prompts = []
472         self.combined_mask = None
473         self.clear_prompts_btn.setEnabled(False)
474         self.save_mask_btn.setEnabled(False)
475         self.toggle_view_btn.setEnabled(False)
476         self.update_canvas()
477         self.status_label.setText("Prompts_cleared")
478
479     def clear_all(self):
480         """Clear everything including the loaded image"""
481         self.original_image = None
482         self.display_image = None
483         self.original_size = None
484         self.prompts = []
485         self.combined_mask = None
486         self.temp_box = None
487
488         # Reset UI

```

```
488         self.canvas.clear()
489         self.canvas.setPixmap(QPixmap())
490         self.run_btn.setEnabled(False)
491         self.clear_prompts_btn.setEnabled(False)
492         self.clear_all_btn.setEnabled(False)
493         self.save_mask_btn.setEnabled(False)
494         self.toggle_view_btn.setEnabled(False)
495         self.status_label.setText("Ready")
496
497     if __name__ == "__main__":
498         app = QApplication(sys.argv)
499         window = SegmentationUI()
500         window.show()
501         sys.exit(app.exec_())
```