

CV Coursework on Image Segmentation

Anonymous Authors *
University of Edinburgh

Abstract

This project investigates several deep learning approaches for semantic image segmentation using the Oxford-IIIT Pet Dataset. We implement and compare UNet-based models, a self-supervised Autoencoder, CLIP feature-driven segmentation, and an interactive prompt-based model. Our experiments evaluate the impact of data augmentation, loss functions, class weighting, and pretraining strategies. Results show that CLIP-based models significantly outperform others by utilizing strong visual representations learned from contrastive pretraining. Data augmentation improves resilience to input perturbations, while prompt-guided segmentation enables effective user interaction. Our findings highlight the value of combining pre-trained features with tailored model design for accurate and robust segmentation.

1. Introduction

Semantic image segmentation is a fundamental task in computer vision that involves classifying each pixel in an image into predefined categories. It plays a crucial role in applications ranging from medical imaging to autonomous navigation and human-computer interaction.

In this work, we address the challenge of segmenting cats and dogs using the Oxford-IIIT Pet Dataset. For this purpose, we explore multiple neural network architectures: a baseline UNet trained from scratch, an Autoencoder whose encoder is reused for downstream segmentation, a hybrid approach combining frozen CLIP visual embeddings with custom decoders, and a prompt-guided interactive segmentation approach. We aim to evaluate these methods in terms of accuracy and robustness to perturbations through multiple metrics.

2. Background

Semantic image segmentation assigns a class label (e.g., 'cat', 'dog', 'background') to every pixel in an image. Deep learning is a primary tool for this, particularly Convolutional

Neural Networks (CNNs). CNNs use convolution operations with learnable filters to automatically detect spatial patterns like edges and textures, making them effective for image data.

The U-Net architecture, designed for segmentation, is frequently used. It employs an encoder to learn contextually rich embeddings and a decoder to produce the final segmentation map. U-Net's key innovation is skip connections, which pass feature maps directly from the encoder to corresponding decoder layers. This allows the reuse of fine grained spatial details, potentially lost during downsampling, combining them with high level semantic context for more accurate segmentation [4].

To further enhance architectures like U-Net, especially in settings with limited labeled data, transfer learning has become an increasingly valuable approach. In this paradigm, models pre-trained on large datasets are reused for related tasks, eliminating the need to train from scratch. Pre-trained networks can serve as feature extractors or be fine-tuned for specific applications, transferring knowledge from high-resource settings to low-resource ones. In computer vision, large vision-language models such as CLIP (Contrastive Language Image Pre-training) [3], provide powerful visual representations learned by aligning images with text, making them well-suited for transfer learning. CLIP's robust image encoder serves as an effective feature extractor, often boosting performance on downstream tasks like segmentation, especially with limited data.

Training deep networks involves iteratively adjusting weights to minimize a loss function that quantifies prediction error against the ground truth; common examples are Mean Squared Error (MSE) and Cross Entropy (CE). Optimization algorithms use gradients, indicating how loss changes with weights, to systematically update weights and reduce error, enabling the network to learn.

3. Methodology

3.1. Hardware and Software Setup

Experiments were conducted using both Kaggle Kernels (with NVIDIA Tesla P100 GPUs) and Google Colab environments (with NVIDIA A100 GPUs and 40GB of VRAM). Development was primarily done on macOS

*B269422, B270876

systems running Python 3.10.13. The implementation was based on PyTorch, with additional libraries including torchvision for model components, NumPy for numerical operations, Pillow for image handling, Hugging Face Transformers to access pre-trained CLIP models, and Matplotlib graph creation.

3.2. Dataset Preprocessing and Data augmentation

We utilized the Oxford-IIIT Pet Dataset [2], which provides cat and dog images paired with pixel level segmentation masks. A 80/20 training-validation split was performed, stratified by breed to ensure similar class representation in both sets. All image-label pairs in the training set were resized to 256x256 pixels, with zero padding to maintain the original aspect ratio. To resize segmentation masks, we consistently use nearest-neighbour interpolation, while bilinear interpolation is applied to the images. In contrast, the validation and test sets are left unaltered.

To increase volume of training data and improve model robustness, data augmentation techniques were applied using the imgaug library, expanding the training set to 6353 pairs. The specific augmentations and their parameters are detailed in Table 1 (with visual examples in Figures 9, 10 in the Appendix). These enhancements were strategically applied to balance the initial 2:1 dog-to-cat ratio in the training data towards a 1:1 ratio.

Finally, for the prompt-based segmentation model, specialized data was generated from the augmented training set. For each image-label pair, two image-heatmap-label triplets were created. The heatmaps simulate point prompts, with the two points per image deliberately placed on different semantic classes (e.g., one on background, one on the animal).

For the prompt-segmentation model, an additional pre-processing step is required: heatmap creation. For each image in the training data, two random points are selected such that they belong to different classes. Corresponding heatmaps are then generated for these points, resulting in two distinct label maps per image. The same procedure is applied to the validation and test sets.

3.3. Training and Evaluation Setup

Data Handling: We implemented two custom PyTorch Dataset classes. The first, for standard segmentation, reads image-label pairs, applying optional augmentations and scaling image pixel values from [0, 255] to [0, 1]. The second, for prompt-based segmentation, reads image-heatmap-label triplets, similarly scaling image and heatmap values. Data is batched using standard DataLoaders. A custom `collate_fn` handles variable image sizes in validation/test sets by batching samples into lists.

Training Loop: Each epoch iterates through the training

DataLoader. The images and labels can optionally be resized with zero padding to meet the model’s requirements. The loop executes a forward pass, computes the loss between the prediction and label, and performs a backward pass to calculate gradients. Gradient accumulation can be adopted, by summing the gradients over several steps before an optimizer update, effectively simulating larger batch sizes. All the provided models are trained for 100 epochs.

Evaluation Loop: After each training epoch, the model is evaluated on the validation set. Input images are resized (bilinear interpolation, with padding) to the model’s expected size, storing original size metadata. The model generates predictions (logits), which are then reverted to their original dimensions using the stored metadata (removing padding, applying bilinear interpolation). These original-sized predictions are compared against the original ground-truth labels.

Metrics: A `MetricsHistory` class accumulates True Positives (TP), False Positives (FP), False Negatives (FN), and True Negatives (TN) per class over the validation set. At the end of an epoch, it computes mean IoU, mean Dice score, and pixel accuracy based on these counts, assessing performance at the original image scale.

Losses: Several loss functions were employed:

- **MSE:** Used for reconstruction. It measures the average squared distance between predicted pixel value, \hat{y} , and true pixel value, y , across the image.

$$\mathcal{L}_{MSE} = \frac{1}{N} \sum_{i \in \text{pixels}} (y_i - \hat{y}_i)^2$$

- **CE:** Used for multi-class segmentation. For a single pixel i with true class c , the model yields logits, \hat{y} . Minimizing the CE loss corresponds to maximizing the predicted probability $P(c|i)$ obtained after applying softmax to these logits.

$$\mathcal{L}_{CE} = -\frac{1}{N} \sum_{i \in \text{pixels}} \log(\text{softmax}(\hat{y}))$$

- **CE+Soft Dice:** To balance pixel-wise accuracy and region overlap, especially under class imbalance, CE loss was combined with Soft Dice loss. Soft Dice loss computes the Dice coefficient directly on predicted probabilities ensuring differentiability. For a single class c : $\mathcal{L}_{Dice,c} = -\frac{2 \sum_{i \in \text{pixels}} P(c|i) y_{ic} + \epsilon}{\sum_{i \in \text{pixels}} P(c|i) + \sum_{i \in \text{pixels}} y_{ic} + \epsilon}$ where y_{ic} is 1 if pixel i belongs to class c and 0 otherwise, and ϵ is a small smoothing constant. This is averaged over the selected classes and added to \mathcal{L}_{CE} .

Table 1. Reasoning and Implementation Details of augmentations.

Augmentation	Implementation Details
Rotation	Decrease dependence on animal position. Rotate images and their labels by a random angle between 45 and 315 degrees, applying zero padding. Some parts of the image, including parts of the object, may be cut. Then resize with zero padding.
Random Masking	Increase robustness by simulating partial occlusions and forcing the model to learn more context-aware representations. Increase ability to segment animals partially covered by objects. Randomly mask (drop) pixels (set to black), on both the image and label. Then resize with zero padding.
Cropping	Learn smaller animal features and parts and provide unpadded training images as input to the model. Extract a square from a random position in the image, where the side of the square is $\frac{2}{3}$ of the shorter side of the original image. Then resize with zero padding.
Resizing+Cropping	Learn smaller animal features and parts and provide unpadded training images as input to the model. Scaling the shorter side to the target size 512 while preserving aspect ratio, followed by center cropping. No additional resize is needed.
Random color jitter	Decrease reliance on color and increase robustness to noise. Randomly adjust each pixel’s color by adding noise sampled from laplace distributions elementwise to the images. Then resize with zero padding.
Grayscale	Decrease reliance on color, learn from texture and shape. Convert images to grayscale.
Blurring	Increase robustness to image quality degradation. Learn global shape and structure. Blur an image by computing means over neighbourhoods.
Contrast Decrease	Decrease reliance on irrelevant fine details and focus on semantically meaningful structures, by making important features stand out under subdued visual conditions. Adjust contrast by scaling each pixel (v) to $127 + \alpha \cdot (v - 127)$.
Merging two images	Concatenate two images side by side, then resize with zero padding. Approximately one-third of the augmented pairs consist of two cats, another third pair a cat with a dog, and the remaining third consist of two dogs. Increase ability to distinguish cats from dogs in the same picture.

$$\mathcal{L}_{CE+Dice} = \mathcal{L}_{CE} + \mathcal{L}_{Dice}$$

- **Class Weighting:** To counteract class imbalance, we experimented with weighted versions of these losses. Instead of a simple average, a weighted average can assign higher importance to less frequent classes. We tested two schemes:

- **Full Weight:** Class weights calculated as inversely proportional to the number of pixels belonging to each class in the training set.
- **Min Weight:** Same as "Full Weight", but the weight for the boundary class (which is ignored during evaluation) was manually set to the minimum weight assigned among the other classes.

Optimizer: We used the AdamW optimizer [1]. AdamW modifies the Adam optimizer by implementing decoupled weight decay. Instead of adding the L2 regularization term to the gradient itself before the adaptive moment updates (as in standard Adam), AdamW applies the weight decay directly to the weights after the gradient-based update step. This distinction often leads to improved generalization performance and more effective regularization compared to standard Adam with L2 regularization.

3.4. UNet-based end-to-end segmentation neural network

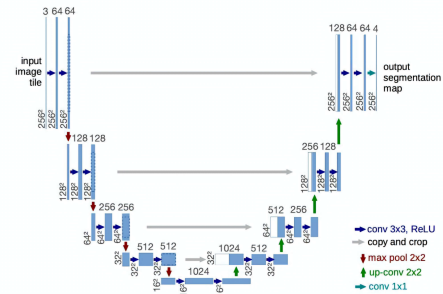


Figure 1. Our U-Net architecture. Each blue box represents a multi-channel feature map, with the number of channels indicated above the box and the spatial dimensions shown at its lower left corner. White boxes denote copied feature maps. Arrows indicate operations in the network. Figure adapted from "U-Net: Convolutional Networks for Biomedical Image Segmentation" [4], using our channels and sizes.

Our U-Net implementation (Figure 1) features a stan-

standard encoder-decoder structure with skip connections. The encoder captures context through five blocks. Each block applies two sequential 3x3 convolutions, each followed by Batch Normalization (BN) and ReLU, using padding to maintain spatial dimensions. The first convolution in each block doubles the number of input feature channels, while the second maintains this channel count. A 2x2 max pooling layer (stride 2) then halves the spatial resolution. This process culminates in a 16x16 bottleneck layer with 1024 channels, holding the most abstract image representation.

The decoder symmetrically reconstructs the segmentation map via four blocks. Each block begins with a 2x2 transposed convolution (stride 2), which doubles spatial resolution and halves the number of input channels. Crucially, features from the corresponding encoder level are concatenated via skip connections. Consistent padding and stride choices ensure spatial alignment for this concatenation, merging high-resolution encoder details with upsampled decoder context. Following concatenation, two sequential 3x3 convolutions (Padding+BN+ReLU) process the merged features; the first of these convolutions halves the channel count (relative to the concatenated feature map), while the second maintains it.

After the final decoder block restores the original 256x256 resolution, a concluding 1x1 convolution maps the 64 feature channels to 4 output channels representing the scores for background, cat, dog, and boundary classes.

For training, we explored several loss functions: standard Cross Entropy, Weighted Cross Entropy (addressing class imbalance), and combinations of CE with Soft Dice Loss (standard or weighted). We employed the AdamW optimizer with a learning rate of 0.001 and weight decay of 0.01. The target effective batch size was 64, achieved using gradient accumulation across multiple steps when necessary due to hardware constraints.

3.5. Autoencoder pre-training for segmentation

As an unsupervised representation learning task, we first train an autoencoder on raw, unlabeled images for image reconstruction. Once trained, the encoder is reused as a frozen feature extractor for the downstream semantic segmentation task with labeled masks. The goal is to learn meaningful latent features and evaluate their effectiveness in transfer learning for segmentation. Training and evaluation were performed using the dataset of augmented 256 × 256 RGB images.

Figure 2 illustrates the models used in the two stages of our complete pipeline:

Stage 1 – Reconstruction Autoencoder (AE):

The reconstruction network follows a symmetric encoder-decoder design similar to U-Net, encouraging the bottleneck to learn latent embeddings from which the

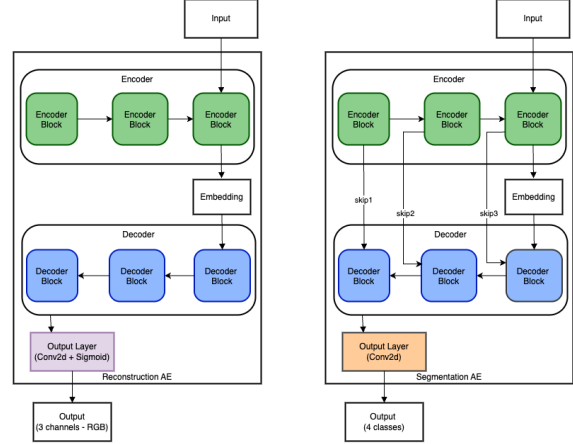


Figure 2. Architecture of Reconstruction AE and Segmentation AE models

entire image can be reconstructed. For reconstruction, we intentionally exclude skip connections. Allowing them would give the model direct access to low-level features, reducing the need for meaningful compression, and potentially leading to less effective latent representations.

Encoder:

The encoder consists of three sequential blocks. Each block contains two convolutional (Conv) layers (kernel size 3, padding 1), each followed by BN and ReLU. A `MaxPool2d` layer at the end of each block halves the spatial resolution.

Starting from an input of shape $(N, 3, 256, 256)$, Block 1 increases channels to 64 and reduces spatial size to 128×128 . The following 2 encoder blocks double the channels and halve the image size again. This leads to a final embedding (bottleneck) of shape $(N, 256, 32, 32)$.

Decoder:

The decoder mirrors the encoder with three upsampling blocks. Each decoder block (`DecoderBlockNoSkips`) uses transposed convolution (`ConvTranspose2d`) with kernel size = 2 and stride = 2 to double height and width. In the first two blocks, channel depth is halved: $256 \rightarrow 128$, then $128 \rightarrow 64$; and in the third block, it is maintained at 64. Each upsampled feature map is refined using two convolutional layers (kernel size 3, padding 1) + BN + ReLU. Finally, a separate `Conv2D` layer maps from 64 channels back to RGB space (3 channels), followed by Sigmoid activation for normalized output: $(N, 64, 256, 256) \rightarrow (N, 3, 256, 256)$.

Reconstruction Training Details:

Reconstruction training was performed using batch size 64, Adam optimizer with learning rate of 0.001 and MSE loss (to promote pixel-wise accuracy). As an additional

experiment, reconstruction was also trained using skip connections. To enable this, the decoder architecture used was identical to that of the segmentation model, employing the same `DecoderBlockWithSkips` structure described in the following section.

Stage 2 – Segmentation Network Using Pre-trained Encoder:

The segmentation model reuses the encoder from the reconstruction AE but adds skip connections in a U-Net-style architecture. These connections merge high-resolution encoder features with upsampled decoder outputs, helping to recover spatial details lost during downsampling. This is essential for semantic segmentation, which relies on precise pixel-level localization, especially along object boundaries.

Encoder:

We reuse the encoder from the reconstruction model by wrapping it in a reusable module (`SegmentationEncoder`) that loads the pre-trained weights from the reconstruction encoder and freezes them during segmentation training.

Decoder with Skip Connections:

Each decoder block (`DecoderBlockWithSkips`) performs upsampling via transposed convolution (kernel size=2, stride=2), doubling spatial dimensions and reducing channels, followed by concatenation with the corresponding skip connection from the encoder and refinement through two stacked Conv–BN–ReLU layers. For example, in the first block: $(N, 256, 32, 32) \xrightarrow{\text{upsample}} (N, 128, 64, 64)$; $\text{Concat}((N, 128, 64, 64), (N, 256, 64, 64)_{\text{skip}_3}) \rightarrow \text{Conv-BN-ReLU} \times 2 \rightarrow (N, 128, 64, 64)$

After 3 decoder blocks, a 1×1 convolution maps the 64 channels to the 4 classes: background (#0), cat (#1), dog (#2), and boundary (#3). Output shape becomes $(N, 4, 256, 256)$, representing class logits per pixel.

Segmentation Training Details:

Segmentation training used a batch size of 64, the AdamW optimizer, with learning rate of 0.001 and weight decay of 0.01, and weighted CE+Soft Dice loss. Both the *Full Weight* and *Min Weight* CE weighting schemes, as defined previously, were compared.

3.6. CLIP features for segmentation

The model’s architecture is depicted in Figure 3. We utilized features from the pre-trained openai/clip-vit-base-patch16 CLIP model, accessed via Hugging Face’s transformers. This Vision Transformer (ViT) was chosen for its strong representations and accessible hidden states. The model processes the input image into 16×16 patches; this patch size, combined with the ViT architecture, results in

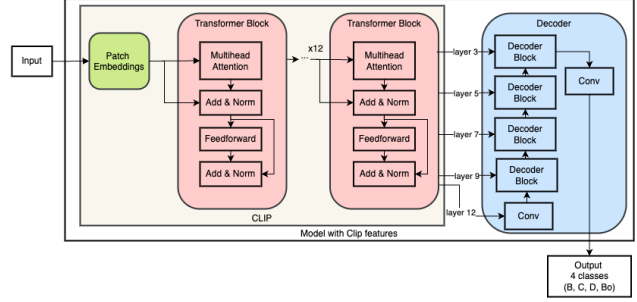


Figure 3. Architecture of Model using CLIP features.

final hidden state outputs (excluding the CLS token) that correspond to 49 patch embeddings (shape 49, 768). These can be reshaped into a 7×7 spatial grid (768, 7, 7) serving as the bottleneck feature map. Crucially, starting from this 7×7 bottleneck dimension allows a UNet-style decoder with four $2 \times$ upsampling stages, as described below, to reconstruct a feature map matching the CLIP encoder’s input resolution (224x224).

Consistent with CLIP requirements, input images were resized to 224x224. The CLIP image encoder’s weights remained frozen, serving purely as a feature extractor. Its final hidden state output (excluding the CLS token) was processed as described above into the (768, 7, 7) spatial grid for the decoder.

This decoder was inspired by the UNet’s expansive path. An initial 1×1 convolution adapted the (768, 7, 7) bottleneck features to 1024 channels. The primary decoder configuration included four up-sampling blocks. Each block performed a 2×2 transposed convolution, followed by concatenation with features extracted from intermediate layers [3, 5, 7, 9] of the frozen CLIP encoder. These intermediate CLIP features (CLS token removed, reshaped, and resized) served as skip connections, aiming to incorporate hierarchical information. A standard double convolution block (3×3 Conv + BN + ReLU, repeated twice) processed the concatenated features. An alternative decoder variant omitted these skip connections, adjusting internal filter counts accordingly. Both decoder versions concluded with a 1×1 convolution projecting features to the 4 required output classes (background, cat, dog, boundary). The overall architecture is depicted in Figure 3.

We employed the training strategy proven effective for the UNet model. Specifically, the weighted version of the combined Cross-Entropy and Dice loss (CE+Dice) was used. Optimization was performed using AdamW with a learning rate of 0.001 and weight decay of 0.01. The model was trained for 100 epochs with a target batch size of 64, employing gradient accumulation as needed due to hardware constraints.

3.7. Prompt-based segmentation

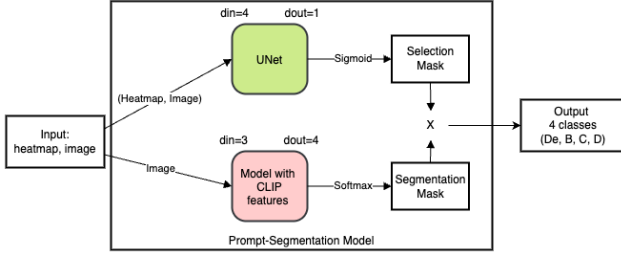


Figure 4. Architecture of Prompt-based Model. Class De - 'deactivated' pixels, B - background/boundary, C - cat, D - dog.

For prompt-based segmentation, we generated triplets of (image, heatmap, label) from the augmented dataset, as stated in 3.2. The heatmap represents the prompt, constructed as a discretized Gaussian ($\sigma=3$) centered on a selected pixel. Augmented images created by merging distinct originals were excluded, as generating accurate ground truth masks for prompts in these complex merged scenarios would require an already perfected segmenter.

The task objective is to segment only the class indicated by the prompt point, treating all other pixels as 'deactivated'. To achieve this, target labels were modified: original background and boundary classes were merged into a single 'background/boundary' class, and a new 'deactivated' class was introduced. Prompts were placed on either the merged background/boundary, cat, or dog regions. While the effective ground truth mask for a given sample highlights only the single selected class against the 'deactivated' background, the model was trained to output scores for all potential semantic classes alongside the 'deactivated' state to retain context. The final label values used for training were: 0 ('deactivated'), 1 ('background/boundary'), 2 ('cat'), and 3 ('dog'). This four-value representation was chosen over a simpler binary output to leverage the detailed cat/dog distinctions for potentially improving shape identification and segmentation accuracy.

We adapted the best-performing non-interactive model, the CLIP feature-based architecture, for this task. The prompt model takes both an image and its corresponding heatmap as input. The pre-existing CLIP-based segmentation component processes the image, outputting initial class segmentation probabilities via Softmax. In parallel, the image is concatenated with the heatmap (4 input channels: R,G,B,Heatmap) and fed into a new 'Selection Network'. This network uses the same UNet decoder architecture but takes 4 input channels and produces a single output channel representing pixel selection probability (via Sigmoid). The architecture of the model is represented in Figure 4.

The final output is computed by element-wise multiplying the class probabilities from the segmentation

component with the selection probability from the Selection Network (broadcast across class channels). This yields a tensor representing joint probabilities $P(\text{class}, \text{selected} | \text{pixel}, \text{prompt})$ for the active classes (background/boundary, cat, dog). The probability for the 'deactivated' class is implicitly one minus the sum of the probability vector. For loss computation, we construct a 4-value probability vector per pixel: $(P(\text{deactivated}), P(\text{bg/boundary, selected}), P(\text{cat, selected}), P(\text{dog, selected}))$.

The network was trained using a combined Negative Log Likelihood (NLL) and Dice Loss applied to this final 4-value probability vector. NLL loss (similar to Cross Entropy, minimizing the negative log probability of the correct class configuration) is suitable for the probabilistic output, while Dice loss optimizes segmentation overlap modulated by the selection. We tested two strategies for the incorporated CLIP component: keeping its weights frozen versus allowing fine-tuning to adapt specifically to the prompt-guided task. Recognizing that correct selection is as critical as classification, equal class weights were applied across background/boundary, cat, and dog classes during loss calculation. Optimization used AdamW (learning rate 0.001, weight decay 0.01) for 100 epochs with an effective batch size of 64 (using gradient accumulation).

The user interface shown in Figure 5 is a locally hosted interactive web application designed to showcase the prompt-based model, allowing users to upload an image, select a point prompt, and view the resulting segmentation mask.

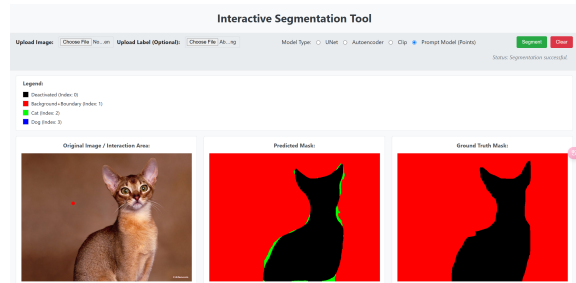


Figure 5. Reconstruction AE and Segmentation AE

4. Evaluation and Discussion of Results

We compared UNet input resolutions of 256x256 and 512x512 to balance accuracy and efficiency. While 512x512 showed marginal gains in Dice/IoU (Figure 8a), it significantly increased computational cost (approx. 4x slower, higher memory). This minor gain is likely because the segmentation task, focused on large coherent regions (animals, background), doesn't heavily rely on fine-grained details. The spatial continuity of classes, combined with in-

terpolation during evaluation, means 256x256 provides sufficient detail. Due to this unfavorable cost-benefit trade-off, we standardized on the efficient 256x256 resolution for subsequent work.

We compared standard Cross Entropy (CE) loss with a combined CE + Soft Dice loss. As shown in Figure 8b, the combined loss yielded slightly superior and more stable metrics (Dice, IoU, Accuracy). We attribute this to CE optimizing pixel-level accuracy while Dice directly encourages better spatial overlap and region consistency, key for segmentation. Dice’s inherent robustness to class imbalance might also contribute to stability. Offering these benefits at no added computational cost, we adopted the combined CE + Soft Dice loss for subsequent work.

To address class imbalance from dominant background pixels, we applied class weighting (Subsection 3.3) to the CE+Soft Dice loss. We tested ‘Full Weight’ (inversely proportional to pixel counts) and ‘Min Weight’ (assigning the minimum weight to the boundary class). Figure 8c shows ‘Full Weight’ performed noticeably better, indicating that assigning higher weights to boundary pixels during training led to improved segmentation performance.

This suggests that accurately learning boundary pixels, though ignored in final evaluation, is crucial. Boundaries define the transition between foreground objects and the background. By precisely learning these boundary regions, the model can better delineate the enclosed objects (cats/dogs) and distinguish them from their surroundings, leading to sharper segmentations and higher IoU/Dice scores. De-emphasizing boundaries (‘Min Weight’) likely permits less precise edge learning, hindering overall segmentation.

Finally, to evaluate the impact of data augmentation on the UNet model under optimal settings (256x256 input, CE+Soft Dice loss, ‘Full Weight’ class weights), we compared its performance with and without augmentation. As shown in Figure 8d, improvements in standard validation metrics (Dice, IoU, Accuracy) were marginal. This unexpected result may be due to the validation set containing clean images, unlike the perturbed ones seen during training. We hypothesize that augmentation primarily enhances model robustness to distortions rather than significantly improving performance on unaltered validation data, a topic explored further in the following subsection.

We evaluated using an Autoencoder (AE) for self-supervised pre-training of the segmentation encoder. Two AE designs were compared during pre-training for image reconstruction: one standard AE (‘AE without Skips’) and one with internal skip connections linking encoder and decoder layers (‘AE with Skips’).

While the ‘AE with Skips’ achieved significantly better reconstruction (lower validation loss, Figure 8e), its utility for the downstream segmentation task proved inferior. After

freezing the pre-trained encoders and training segmentation decoders, the encoder pre-trained without skip connections yielded markedly better segmentation performance (Dice, IoU, Accuracy) than the one pre-trained with skips (Figure 8f).

This inverse relationship suggests the skip connections in the AE, although improving reconstruction fidelity, likely acted as information shortcuts. By allowing the decoder to directly access low-level features from the encoder, these connections reduced the burden on the encoder to learn a compressed, semantically rich representation in the bottleneck. In contrast, training without skip connections forced the encoder to rely solely on its bottleneck for reconstruction, resulting in more robust and transferable features that were better suited for segmentation, despite a modest decrease in reconstruction quality.

For CLIP-based models leveraging a pre-trained encoder, we first examined the impact of adding skip connections between frozen CLIP layers and the custom UNet-style decoder. As shown in Figure 7, the variant with skip connections (using features from CLIP layers [3, 5, 7, 9]) consistently outperformed the version without them across all metrics (Dice, IoU, Accuracy).

This highlights the value of multi-level feature fusion in segmentation, as combining high-level semantic features with intermediate spatial details, similar to standard UNet architecture, proves beneficial. Beyond this architectural similarity, they also help overcome a key weakness of vision transformers. Because CLIP encoders work with fixed-size image patches, discarding precise pixel-wise alignment early in the encoding process, and rely on global attention, they tend to lose fine spatial detail. By bringing in features from earlier layers, where more local structure is preserved, the decoder can recover important details needed for accurate pixel-level predictions.

In the same experiment (Figure 7), we also re-evaluated the impact of class weighting schemes (‘Full Weight’ vs. ‘Min Weight’). Unlike with UNet, where ‘Full Weight’ clearly performed better, the CLIP-based model showed negligible performance differences between the two (e.g., solid vs. dashed blue/red lines), suggesting that CLIP’s rich pre-trained features may already distinguish classes effectively, reducing sensitivity to loss weighting.

Finally, we assessed the effect of data augmentation on the best-performing CLIP-based model configuration (with skip connections and ‘Full Weight’). As shown in Figure 8g, augmentation yielded only minor improvements in standard validation metrics, consistent with UNet results, supporting the hypothesis that its primary benefit lies in improving robustness to perturbations rather than boosting performance on clean validation data.

Since the prompt-based model builds on the CLIP-based segmentation architecture, we evaluated how its training

Table 2. Comparison of Model Performance Metrics on Test Set

Model	Accuracy	Dice	IoU
UNET (Aug)	0.9462	0.8661	0.7687
UNET (No Aug)	0.9444	0.8632	0.7643
CLIP (Aug)	0.9732	0.9442	0.8946
CLIP (No Aug)	0.9723	0.9414	0.8897
Autoencoder	0.8712	0.6804	0.5382
Prompt	0.8321	0.7088	0.5497

state affects final performance. Specifically, we compared keeping the CLIP component frozen (using weights from its initial non-interactive training) versus fine-tuning it alongside the newly introduced Selection Network.

As shown in Figure 8h, fine-tuning led to consistently better performance across all metrics. While the frozen CLIP component was already effective, allowing it to adapt within the context of prompt-based training improved both segmentation accuracy and coordination with the Selection Network. This highlights that joint optimization enhances synergy between components, making fine-tuning essential for optimal interactive segmentation. Fine-tuning adapts CLIP’s global semantic representations toward the spatially localized task of prompt-guided segmentation, where precise alignment between image features and user-specified points is critical. Therefore, the best performing Prompt-based model utilized fine-tuning for the integrated CLIP component.

After identifying optimal configurations, we now summarize their final performance on the held out test set. Table 2 reports key metrics, Accuracy, Dice score, and IoU, for the best-performing variant of each model: UNet (with/without augmentation), CLIP-based model (with/without augmentation), Autoencoder and Prompt-based model.

The performance comparison presented in Table 2 clearly indicates the superiority of the CLIP-based model, which significantly outperformed all other approaches across accuracy, Dice, and IoU metrics, achieving a notably high mIoU of approximately 0.89. This strong performance is attributed to the powerful, generalizable image representations learned by the CLIP encoder during its extensive pre-training on a massive and diverse dataset, coupled with its Vision Transformer architecture’s ability to effectively capture global context.

In comparison, the standard UNet model performed respectably, validating its architecture with skip connections for this task and surpassing the Autoencoder approach. However, lacking large-scale pre-training, it could not reach the segmentation quality of the CLIP-based model.

Consistent with earlier experiments, data augmentation provided only marginal benefits for both UNet and CLIP mod-

els when evaluated on this clean test set.

The Autoencoder pre-trained model yielded significantly weaker results, likely because its features were learned solely through reconstruction on the comparatively small Oxford-IIIT dataset using a less complex architecture, resulting in less robust representations than CLIP.

Finally, the Prompt-based model’s metrics appear lower but must be interpreted in the context of its distinct task, which combines selection based on a heatmap prompt with segmentation. Evaluated on a derived dataset focusing only on the prompted class, its scores reflect performance on this more complex, interactive objective and are not directly comparable to the full semantic segmentation models, yet demonstrate its effectiveness in that specific setting.

As evidenced above, the transfer learning approach yielded significantly different results when using the pre-trained CLIP encoder versus the reconstruction autoencoder (AE) for image segmentation, achieving the highest and lowest IoU scores, respectively. A major factor is the scale and diversity of CLIP’s training dataset, compared to the smaller dataset used for training the AE. Additionally, the pretraining tasks differ in relevance: while the AE is trained for image reconstruction, capturing low-level structural features, CLIP is trained for image classification, a task more closely related to classification. This enables CLIP to learn high-level semantic features that are more transferable to segmentation tasks. For example, by learning to classify an image as a “cat,” CLIP develops features useful not only for identification, helping distinguish between the cat and dog classes) but also for delineation in segmentation. This greater task alignment contributes to why the transfer learning was more effective for the CLIP-based model than the AE.

4.1. Robustness exploration

To evaluate generalization and the impact of data augmentation on resilience, we subjected the CLIP (Aug/No Aug) and UNet (Aug/No Aug) models to eight perturbation types: Gaussian pixel noise, Gaussian blurring, contrast increase/decrease, brightness increase/decrease, occlusion, and salt and pepper noise, each applied at ten increasing severity levels. Performance (Mean Dice) against perturbation level is plotted in Figure 6.

The results consistently demonstrate the value of augmentation: across most perturbations, the augmented models (solid lines) maintained higher performance or degraded less rapidly than their non-augmented counterparts (dashed lines), strongly supporting our hypothesis that augmentation enhances robustness.

Model stability varied by perturbation. Gaussian pixel noise, contrast increase, brightness changes, and occlusion generally caused only minor or gradual declines, with augmented models showing a slight advantage. However, per-

Model Robustness Comparison: Dice Score vs Perturbation Level

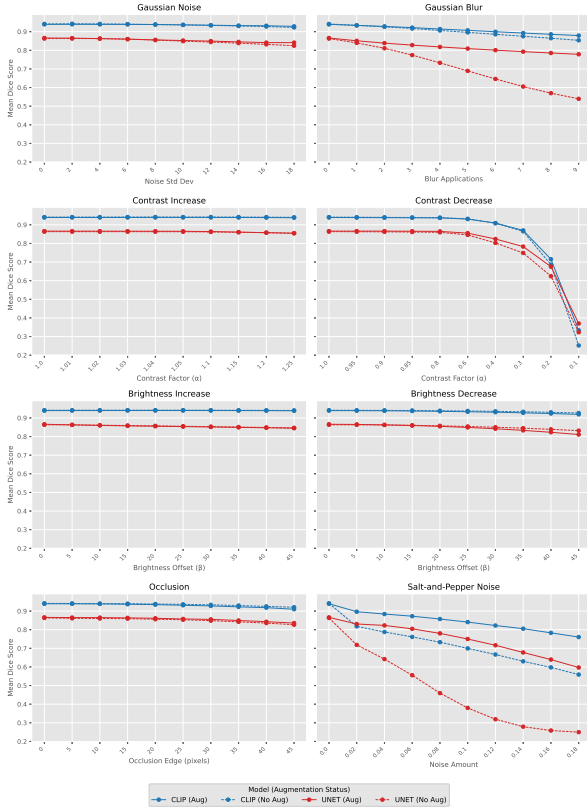


Figure 6. Comparison of CLIP-based model and UNet, with and without augmentation, under increasing perturbation levels.

formance dropped significantly under Gaussian blur and salt and pepper noise, particularly for non-augmented models (especially UNet). Notably, under these specific strong perturbations, the augmented UNet occasionally surpassed the non-augmented CLIP model, highlighting how targeted augmentation can provide substantial resilience against certain degradations, sometimes outweighing general pre-training benefits.

Universally, contrast decrease posed the greatest challenge. All models experienced sharp performance plunges as contrast approached zero. This underscores the fundamental reliance of segmentation algorithms on detecting discernible edges, which are defined by differences in pixel intensity (i.e., contrast). Severely reducing contrast effectively diminishes or erases this critical edge information, crippling the models’ ability to accurately delineate object boundaries, regardless of the training strategy employed.

In summary, these experiments quantitatively confirm that data augmentation significantly boosts model robustness, especially against spatial noise like blur and salt and pepper. While CLIP’s pre-training provides superior perfor-

Clip + UNet Decoder: Skip Connections & Dice Weighting Comparisons

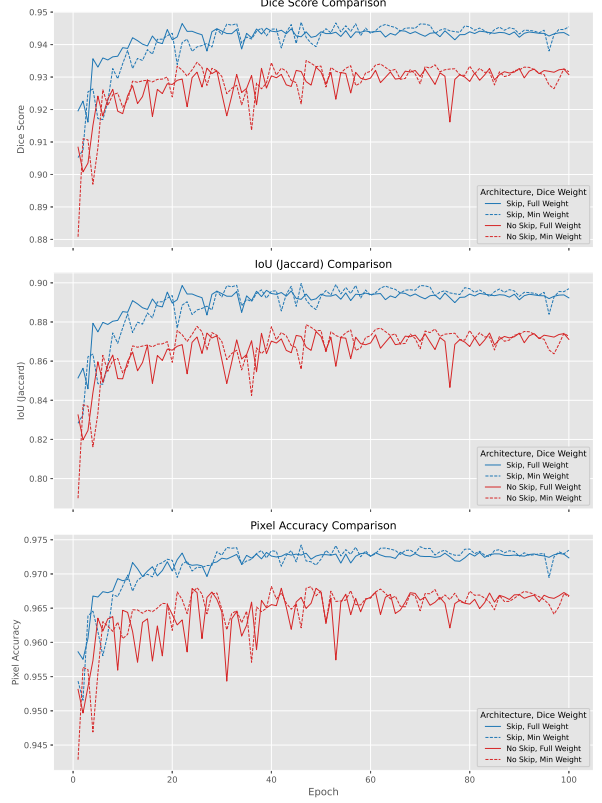
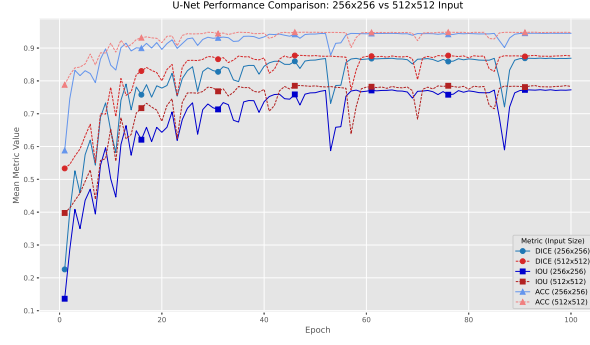


Figure 7. Comparison of the CLIP-based model performance with vs. without skip connections, using CE+Soft Dice loss with simple class weights vs. reduced boundary weight.

mance on clean data, augmentation adds a crucial layer of resilience. Severe contrast reduction remains a fundamental challenge for these segmentation approaches.

5. Conclusion

This work explored deep learning approaches for image segmentation, showing that models using pre-trained CLIP encoders outperform others by transferring semantically rich representations learned via large-scale vision-language classification. Skip connections in transformer-based decoders recovered spatial detail, and retaining boundary contributions via class weighting enhances boundary precision. While the UNet remained competitive, its reliance on lower-level features made it less adaptable than CLIP-based models. Fine-tuning pre-trained components significantly improved prompt-guided segmentation, where spatial reasoning is vital, and augmentation significantly enhanced robustness to noise and distributional shifts. Overall, our results show the value of combining pre-trained visual features with well-designed architectures and training for robust and accurate segmentation.



(a) Comparison of UNet performance metrics when trained on 256x256 vs. 512x512 input images.



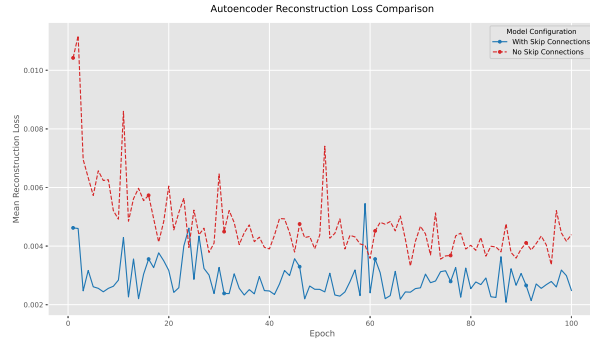
(b) Comparison of UNet performance using CE loss vs. CE + Soft Dice loss.



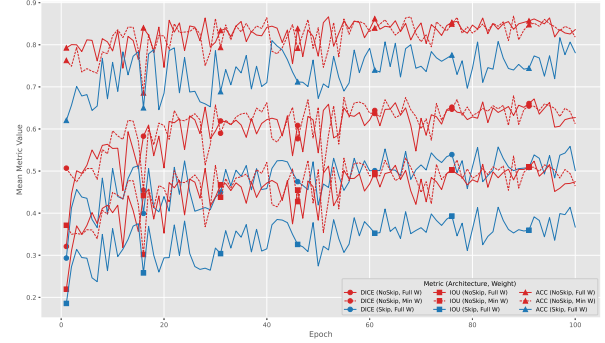
(c) Comparison of UNet performance with CE+Soft Dice loss using simple class weights vs. reduced boundary weight.



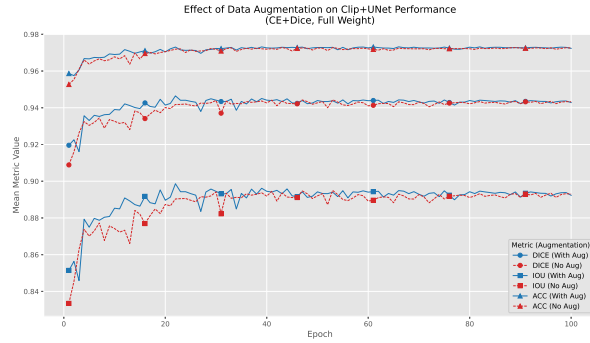
(d) Comparison between the UNet performance with vs. without data augmentation.



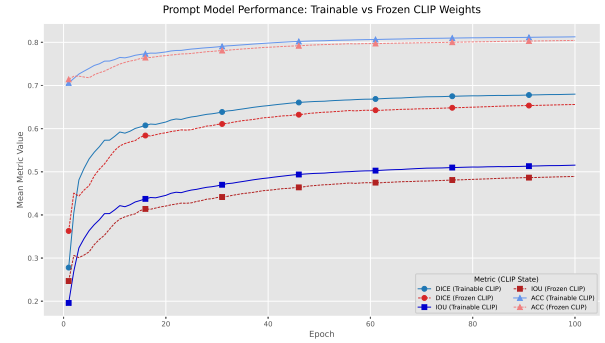
(e) Comparison of AE Reconstruction performance (validation loss) with vs without skip connections.



(f) Comparison of AE Segmentation performance when using an Encoder Pre-trained with vs without skip connections.



(g) CLIP-based model performance with vs. without data augmentation.



(h) Prompt-based model performance with frozen vs. fine-tuned CLIP weights.

Figure 8. Performance Experiments for UNet (rows 1-2), AE (3rd row), CLIP-based (bottom left) and Prompt-based (bottom right) models.

6. Appendix

Division of Responsibilities

The work was divided equally between both group members, with each taking responsibility for specific components and contributing jointly to others. Both students collaborated on writing the report and worked together on data augmentation techniques and the development of the UNet model. Specifically, Student B269422 focused on implementing and refining the encoder part of the UNet architecture, while Student B270876 concentrated on optimizing the decoder section.

Student B269422 was solely responsible for designing and implementing the Autoencoder module, as well as developing the user interface (UI) for interacting with the models. On the other hand, Student B270876 took charge of building and training two alternative models: one utilizing CLIP features and another based on prompt-based learning strategies.

Both students independently ran evaluation experiments to assess model performance. They also explored robustness under varying conditions in parallel by dividing up datasets and scenarios to test against. The analysis of these evaluations was later compiled collaboratively into the final report.

Augmentation Images

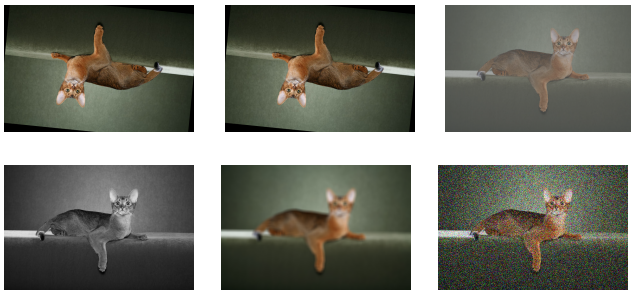


Figure 9. Examples of image augmentations used in training, including rotation, cropping, contrast decrease (top row), and grayscale, blur, and noise/color jitter (bottom row). These are not yet resized to 256×256 .

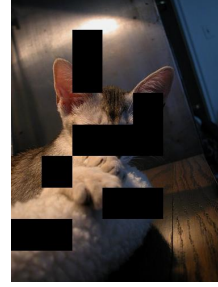


Figure 10. Example of image augmented with random masking, not yet resized to 256×256 .

References

- [1] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization, 2019. [3](#)
- [2] Omkar M. Parkhi, Andrea Vedaldi, Andrew Zisserman, and C. V. Jawahar. Cats and dogs. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2012. [2](#)
- [3] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision, 2021. [1](#)
- [4] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation, 2015. [1, 3](#)

Code

UNet (unet/unet.py)

```

1 import torch
2 from torch import nn
3
4 class DoubleConvReLU(nn.Module):
5     """
6     Applies two convolutional layers with ReLU activation and batch normalization.
7     Args:
8         din (int): Number of input channels.
9         dout (int): Number of output channels.
10    Returns:
11        torch.Tensor: Output tensor after applying the double convolution with ReLU.
12    """
13    def __init__(self, din, dout):
14        super().__init__()
15        self.doubleConvReLU = nn.Sequential(
16            nn.Conv2d(din, dout, kernel_size=3, padding=1),
17            nn.BatchNorm2d(dout),
18            nn.ReLU(),
19            nn.Conv2d(dout, dout, kernel_size=3, padding=1),
20            nn.BatchNorm2d(dout),
21            nn.ReLU(),
22        )
23
24    def forward(self, x):
25        return self.doubleConvReLU(x)
26
27
28 class Down(nn.Module):
29     """
30     Downsamples the input using max pooling and applies double convolution.
31     Args:
32         din (int): Number of input channels.
33         dout (int): Number of output channels.
34    Returns:
35        torch.Tensor: Output tensor after downscaling and double convolution.
36    """
37    def __init__(self, din, dout):
38        super().__init__()
39        self.maxpool_doubleConv = nn.Sequential(
40            nn.MaxPool2d(kernel_size=2, stride=2),
41            DoubleConvReLU(din, dout)
42        )
43
44    def forward(self, x):
45        return self.maxpool_doubleConv(x)
46
47
48 class Up(nn.Module):
49     """
50     Upscales the input using transposed convolution, concatenates with a skip connection, and applies double convolution.
51     Args:
52         din (int): Number of input channels.
53         dout (int): Number of output channels.
54    Returns:
55        torch.Tensor: Output tensor after upscaling, concatenation, and double convolution.
56    """
57    def __init__(self, din, dout):
58        super().__init__()
59
60        self.upsample = nn.ConvTranspose2d(din, dout, kernel_size=2, stride=2)
61        self.doubleConv = DoubleConvReLU(din, dout)
62
63    def forward(self, x1, x2):
64        x = torch.cat([x1, self.upsample(x2)], dim=1)
65        return self.doubleConv(x)
66

```



```

67 class unet(nn.Module):
68     """
69     Implements the U-Net architecture.
70     Args:
71         din (int): Number of input channels.
72         dout (int): Number of output channels.
73     Returns:
74         torch.Tensor: Output tensor after passing through the U-Net.
75     """
76     def __init__(self, din, dout):
77         super().__init__()
78         self.scale = 1
79
80         self.down1 = DoubleConvReLU(din, self.scale * 64)
81         self.down2 = Down(self.scale * 64, self.scale * 128)
82         self.down3 = Down(self.scale * 128, self.scale * 256)
83         self.down4 = Down(self.scale * 256, self.scale * 512)
84         self.down5 = Down(self.scale * 512, self.scale * 1024)
85
86         self.up1 = Up(self.scale * 1024, self.scale * 512)
87         self.up2 = Up(self.scale * 512, self.scale * 256)
88         self.up3 = Up(self.scale * 256, self.scale * 128)
89         self.up4 = Up(self.scale * 128, self.scale * 64)
90
91         self.output = nn.Conv2d(self.scale * 64, dout, kernel_size=1)
92
93     def forward(self, x):
94         x1 = self.down1(x)
95         x2 = self.down2(x1)
96         x3 = self.down3(x2)
97         x4 = self.down4(x3)
98         x5 = self.down5(x4)
99
100         x = self.up1(x4, x5)
101         x = self.up2(x3, x)
102         x = self.up3(x2, x)
103         x = self.up4(x1, x)
104
105         return self.output(x)

```

Autoencoder (autoencoder/autoencoder.py)

```

1 import torch.nn as nn
2 import numpy as np
3 import torch
4
5
6 class EncoderBlock(nn.Module):
7     """
8     A block for the encoder, consisting of two convolutional layers, batch normalization, ReLU activations, and max pooling.
9     Args:
10         din (int): Number of input channels.
11         dout (int): Number of output channels.
12     Returns:
13         tuple: A tuple containing the pooled output and the skip connection features.
14     """
15     def __init__(self, din, dout):
16         super().__init__()
17         self.conv1 = nn.Conv2d(din, dout, kernel_size=3, padding=1, bias=False)
18         self.bn1 = nn.BatchNorm2d(dout)
19         self.relu1 = nn.ReLU()
20         self.conv2 = nn.Conv2d(dout, dout, kernel_size=3, padding=1, bias=False)
21         self.bn2 = nn.BatchNorm2d(dout)
22         self.relu2 = nn.ReLU(inplace=True)
23         self.pool = nn.MaxPool2d(kernel_size=2, stride=2)
24

```

```

25     def forward(self, x):
26         x = self.conv1(x)
27         x = self.bn1(x)
28         x = self.relu1(x)
29         x = self.conv2(x)
30         x = self.bn2(x)
31         skip_connection = self.relu2(x) # Features before pooling
32         pooled_output = self.pool(skip_connection)
33         return pooled_output, skip_connection
34
35     class Encoder(nn.Module):
36         """
37         The encoder module composed of multiple encoder blocks.
38         Args:
39             din (int): Number of input channels.
40             base_channels (int): Base number of channels to be used for the encoder blocks.
41         Returns:
42             tuple: A tuple containing the bottleneck features and skip connections from each encoder block.
43         """
44         def __init__(self, din, base_channels):
45             super().__init__()
46             self.encoderPart1 = EncoderBlock(din, base_channels)
47             self.encoderPart2 = EncoderBlock(base_channels, base_channels*2)
48             self.encoderPart3 = EncoderBlock(base_channels*2, base_channels*4)
49
50         def forward(self, x):
51             x1_pooled, skip1 = self.encoderPart1(x)
52             x2_pooled, skip2 = self.encoderPart2(x1_pooled)
53             bottleneck, skip3 = self.encoderPart3(x2_pooled)
54             return bottleneck, skip3, skip2, skip1
55
56     class DecoderBlockWithSkips(nn.Module):
57         """
58         A decoder block that uses skip connections.
59         Args:
60             din_up (int): Number of input channels from the previous upsampled layer.
61             din_skip (int): Number of input channels from the skip connection.
62             dout (int): Number of output channels.
63         Returns:
64             torch.Tensor: Output tensor after upsampling, concatenation, and convolution.
65         """
66         def __init__(self, din_up, din_skip, dout):
67             super().__init__()
68             self.up = nn.ConvTranspose2d(din_up, dout, kernel_size=2, stride=2)
69             conv_input_channels = dout + din_skip # Input to convs includes skip features
70             self.convs = nn.Sequential(
71                 nn.Conv2d(conv_input_channels, dout, kernel_size=3, padding=1, bias=False),
72                 nn.BatchNorm2d(dout),
73                 nn.ReLU(inplace=True),
74                 nn.Conv2d(dout, dout, kernel_size=3, padding=1, bias=False),
75                 nn.BatchNorm2d(dout),
76                 nn.ReLU(inplace=True)
77             )
78
79         def forward(self, x, skip_features):
80             x_upsampled = self.up(x)
81             if skip_features.shape[2:] != x_upsampled.shape[2:]:
82                 diffY = skip_features.size()[2] - x_upsampled.size()[2]
83                 diffX = skip_features.size()[3] - x_upsampled.size()[3]
84                 if diffY < 0 or diffX < 0:
85                     raise ValueError("Upsampled larger than skip")
86                 skip_features = skip_features[:, :, diffY // 2 : diffY // 2 + x_upsampled.size()[2],
87                                         diffX // 2 : diffX // 2 + x_upsampled.size()[3]]
88
89             x_concat = torch.cat([x_upsampled, skip_features], dim=1)
90
91

```

```

92         output = self.convs(x_concat)
93         return output
94
95
96 class DecoderWithSkips(nn.Module):
97     """
98     The decoder module with skip connections.
99     Args:
100         base_channels (int): Base number of channels.
101     Returns:
102         torch.Tensor: Output tensor after decoding.
103     """
104     def __init__(self, base_channels):
105         super().__init__()
106         self.decoderBlock1 = DecoderBlockWithSkips(din_up=base_channels*4, din_skip=base_channels*4, dout=base_channels*4)
107         self.decoderBlock2 = DecoderBlockWithSkips(din_up=base_channels*2, din_skip=base_channels*2, dout=base_channels*2)
108         self.decoderBlock3 = DecoderBlockWithSkips(din_up=base_channels, din_skip=base_channels, dout=base_channels)
109
110     def forward(self, bottleneck, skip3, skip2, skip1):
111         d1 = self.decoderBlock1(bottleneck, skip3)
112         d2 = self.decoderBlock2(d1, skip2)
113         d3 = self.decoderBlock3(d2, skip1)
114         return d3 # Output feature map (B, base_channels, H, W)
115
116
117 class DecoderBlockNoSkips(nn.Module):
118     """
119     A decoder block without skip connections.
120     Args:
121         din_up (int): Number of input channels from the previous upsampled layer.
122         dout (int): Number of output channels.
123     Returns:
124         torch.Tensor: Output tensor after upsampling and convolution.
125     """
126     def __init__(self, din_up, dout):
127         super().__init__()
128         # Upsample and change channels
129         self.up = nn.ConvTranspose2d(din_up, dout, kernel_size=2, stride=2)
130         # Convolutions only process the upsampled features
131         # Input channels to convs is just 'dout' (output channels of self.up)
132         self.convs = nn.Sequential(
133             nn.Conv2d(dout, dout, kernel_size=3, padding=1, bias=False),
134             nn.BatchNorm2d(dout),
135             nn.ReLU(inplace=True),
136             nn.Conv2d(dout, dout, kernel_size=3, padding=1, bias=False),
137             nn.BatchNorm2d(dout),
138             nn.ReLU(inplace=True)
139         )
140
141     def forward(self, x):
142         # Only takes features from the previous layer 'x'
143         x_upsampled = self.up(x)
144         # No concatenation
145         output = self.convs(x_upsampled)
146         return output
147
148
149 class DecoderNoSkips(nn.Module):
150     """
151     The decoder module without skip connections.
152     Args:
153         base_channels (int): Base number of channels.
154     Returns:
155         torch.Tensor: Output tensor after decoding.
156     """
157     def __init__(self, base_channels):
158         super().__init__()

```

```

159         self.decoderBlock1 = DecoderBlockNoSkips(din_up=base_channels*4, dout=base_channels*2) # 256 -> 128
160         self.decoderBlock2 = DecoderBlockNoSkips(din_up=base_channels*2, dout=base_channels) # 128 -> 64
161         self.decoderBlock3 = DecoderBlockNoSkips(din_up=base_channels, dout=base_channels) # 64 -> 64
162
163     def forward(self, bottleneck):
164         # Only takes the bottleneck as input
165         d1 = self.decoderBlock1(bottleneck)
166         d2 = self.decoderBlock2(d1)
167         d3 = self.decoderBlock3(d2)
168         return d3 # Output feature map (B, base_channels, H, W)
169
170
171 class ReconstructionAutoencoder(nn.Module):
172     """
173     A reconstruction autoencoder model.
174     Args:
175         din (int): Number of input channels.
176         dout (int): Number of output channels.
177         base_channels (int): Base number of channels.
178     Returns:
179         torch.Tensor: Reconstructed output tensor.
180     """
181     def __init__(self, din, dout=3, base_channels=64):
182         super().__init__()
183         self.encoder = Encoder(din, base_channels)
184         # No skips in Reconstruction AE
185         self.decoder = DecoderNoSkips(base_channels)
186
187         # Final layer maps DecoderSimple output (base_channels=64) to reconstruction
188         self.decoderOut = nn.Sequential(
189             nn.Conv2d(base_channels, dout, kernel_size=3, padding=1),
190             nn.Sigmoid()
191         )
192
193     def forward(self, x):
194         # Encode, getting bottleneck and skip connections
195         bottleneck, skip3_ignored, skip2_ignored, skip1_ignored = self.encoder(x)
196         # Decode using ONLY the bottleneck with the simple decoder
197         decoded_features = self.decoder(bottleneck)
198         # Apply final layer
199         reconstructed = self.decoderOut(decoded_features)
200         return reconstructed
201
202
203 class SegmentationEncoder(nn.Module):
204     """
205     An encoder module for segmentation tasks, with optional pre-trained weights and freezing.
206     Args:
207         din (int): Number of input channels.
208         base_channels (int): Base number of channels.
209         pretrained_encoder_path (str, optional): Path to the pre-trained encoder weights. Defaults to None.
210         freeze_encoder (bool, optional): Whether to freeze the encoder parameters. Defaults to True.
211     Returns:
212         tuple: Bottleneck and skip connections from the encoder.
213     """
214     def __init__(self, din, base_channels, pretrained_encoder_path=None, freeze_encoder=True):
215         super().__init__()
216         self.encoder = Encoder(din, base_channels)
217
218         if pretrained_encoder_path:
219             try:
220                 full_state_dict = torch.load(pretrained_encoder_path, weights_only=False, map_location=lambda storage,
221                 # Handle potential checkpoint structure variations
222                 if "model_state_dict" in full_state_dict:
223                     model_state_dict = full_state_dict["model_state_dict"]
224                 elif "state_dict" in full_state_dict:
225                     model_state_dict = full_state_dict["state_dict"]

```



```

226         else:
227             model_state_dict = full_state_dict # Assume it's the state dict directly
228
229         encoder_state_dict = {}
230         has_encoder_prefix = any(k.startswith('encoder.') for k in model_state_dict.keys())
231
232         for key, value in model_state_dict.items():
233             if has_encoder_prefix:
234                 if key.startswith('encoder.'):
235                     new_key = key[len('encoder.'):]
236                     encoder_state_dict[new_key] = value
237
238         if not encoder_state_dict:
239             print("Warning: Could not extract encoder state dict. Checkpoint might be empty or incompatible.")
240         else:
241             load_result = self.encoder.load_state_dict(encoder_state_dict, strict=True) # Use strict=False for
242             print(f"Loaded encoder weights. Load result:")
243             if load_result.missing_keys:
244                 print(" Missing keys:", load_result.missing_keys)
245             if load_result.unexpected_keys:
246                 print(" Unexpected keys:", load_result.unexpected_keys)
247             if not load_result.missing_keys and not load_result.unexpected_keys:
248                 print(" All keys matched successfully.")
249
250         except FileNotFoundError:
251             print(f"Warning: Pre-trained encoder file not found: {pretrained_encoder_path}. Using random weights.")
252         except Exception as e:
253             print(f"Warning: Error loading weights: {e}. Check compatibility. Using random weights.")
254
255         if freeze_encoder:
256             if not pretrained_encoder_path:
257                 print("Warning: Freezing encoder, but no pre-trained weights were loaded.")
258             for param in self.encoder.parameters():
259                 param.requires_grad = False
260             print("Encoder parameters frozen.")
261         else:
262             print("Encoder parameters are trainable.")
263
264
265     def forward(self, x):
266         # Returns bottleneck, skip3, skip2, skip1
267         return self.encoder(x)
268
269
270
271 class SegmentationAutoencoder(nn.Module):
272     """
273     A segmentation autoencoder model.
274     Args:
275         din (int): Number of input channels.
276         base_channels (int): Base number of channels.
277         num_classes (int): Number of output classes.
278         pretrained_encoder_path (str, optional): Path to the pre-trained encoder weights. Defaults to None.
279         freeze_encoder (bool, optional): Whether to freeze the encoder parameters. Defaults to True.
280     Returns:
281         torch.Tensor: Segmentation logits.
282     """
283     def __init__(self, din, base_channels=64, num_classes=4, pretrained_encoder_path=None, freeze_encoder=True):
284         super().__init__()
285         self.num_classes = num_classes
286
287         # Initialize the Encoder (via wrapper for loading/freezing)
288         self.encoder = SegmentationEncoder(din, base_channels, pretrained_encoder_path=pretrained_encoder_path, freeze
289
290         # Use the Decoder WITH Skips for Segmentation
291         self.decoder = DecoderWithSkips(base_channels)
292

```

```

293     # Final convolution maps DecoderWithSkips output (base_channels=64) to class scores
294     self.finalConv = nn.Conv2d(base_channels, num_classes, kernel_size=1)
295
296     def forward(self, x):
297         # 1. Encoder gets bottleneck and skips
298         bottleneck, skip3, skip2, skip1 = self.encoder(x)
299
300         # 2. Decoder uses bottleneck AND skips
301         decoder_output = self.decoder(bottleneck, skip3, skip2, skip1)
302
303         # 3. Final 1x1 convolution for class logits
304         segmentation_logits = self.finalConv(decoder_output)
305
306     return segmentation_logits

```

Clip (clip/clipunet.py)

```

1  import torch
2  import torch.nn as nn
3  import torch.nn.functional as F
4  from transformers import CLIPVisionModel, CLIPVisionConfig
5
6
7  class ClipViTEncoder(nn.Module):
8      """
9      Encodes an image using CLIP's Vision Transformer.
10
11      Args:
12          model_name (str): The name of the pre-trained CLIP model to use.
13          freeze_encoder (bool): Whether to freeze the CLIP encoder's parameters.
14          skip_indices (list): The indices of the hidden states to use for skip connections.
15
16      Returns:
17          bottleneck_features: The bottleneck features from the encoder.
18          skip_features_list: A list of skip features from the encoder.
19      """
20     def __init__(self, model_name="openai/clip-vit-base-patch16", freeze_encoder=True, skip_indices=[3, 5, 7, 9]):
21         super().__init__()
22
23         self.skip_indices = sorted(skip_indices)
24
25         self.config = CLIPVisionConfig.from_pretrained(model_name)
26         self.clip_vit = CLIPVisionModel.from_pretrained(model_name)
27
28         if freeze_encoder:
29             for param in self.clip_vit.parameters():
30                 param.requires_grad = False
31
32         self.grid_size = self.config.image_size // self.config.patch_size
33         self.hidden_dim = self.config.hidden_size
34
35     def forward(self, x):
36         if x.shape[2] != self.config.image_size or x.shape[3] != self.config.image_size:
37             print(
38                 f"Input image size ({x.shape[2]}x{x.shape[3]}) doesn't match "
39                 f"CLIP expected size ({self.config.image_size}x{self.config.image_size}). "
40                 f"Behavior may be unexpected. Consider resizing input."
41             )
42
43         outputs = self.clip_vit(pixel_values=x, output_hidden_states=True)
44         all_hidden_states = outputs.hidden_states
45         last_hidden_state = outputs.last_hidden_state
46
47         # (N, 49, 768) -> (N, 7, 7, 768) -> (N, 768, 7, 7) = (N, C, H, W)
48         patch_embeddings = last_hidden_state[:, 1:, :] # Remove CLS
49         bottleneck_features = patch_embeddings \

```

```

50         .reshape(x.shape[0], self.grid_size, self.grid_size, self.hidden_dim) \
51         .permute(0, 3, 1, 2).contiguous()
52
53     skip_features_list = []
54     for i in self.skip_indices:
55         hidden_state = all_hidden_states[i]
56         patch_embeddings = hidden_state[:, 1:, :] # Remove CLS
57
58         # (N, 49, 768) -> (N, 7, 7, 768) -> (N, 768, 7, 7) = (N, C, H, W)
59         reshaped_features = patch_embeddings \
60             .reshape(x.shape[0], self.grid_size, self.grid_size, self.hidden_dim) \
61             .permute(0, 3, 1, 2).contiguous()
62
63         skip_features_list.append(reshaped_features)
64
65     return bottleneck_features, skip_features_list
66
67
68 class DecoderBlock(nn.Module):
69     """
70     A single decoder block for the UNet.
71
72     Args:
73         in_channels (int): The number of input channels.
74         in_channels_skip (int): The number of channels from the skip connection.
75         out_channels (int): The number of output channels.
76
77     Returns:
78         x (torch.Tensor): The output tensor after the decoder block.
79     """
80     def __init__(self, in_channels, in_channels_skip, out_channels):
81         super().__init__()
82
83         self.upsample = nn.ConvTranspose2d(in_channels, in_channels // 2, kernel_size=2, stride=2)
84         self.skip_conv = nn.Conv2d(in_channels_skip, in_channels // 2, kernel_size=1) # 768 channels to in_channels//2
85
86         self.conv_block = nn.Sequential(
87             nn.Conv2d(in_channels, out_channels, kernel_size=3, padding=1, bias=False),
88             nn.BatchNorm2d(out_channels),
89             nn.ReLU(inplace=True),
90             nn.Conv2d(out_channels, out_channels, kernel_size=3, padding=1, bias=False),
91             nn.BatchNorm2d(out_channels),
92             nn.ReLU(inplace=True),
93         )
94
95     def forward(self, x, skip):
96         x = self.upsample(x)
97         skip = self.skip_conv(skip)
98
99         if skip.shape[2:] != x.shape[2:]:
100             skip= F.interpolate(skip, size=x.shape[2:], mode='bilinear', align_corners=False)
101
102         x = torch.cat([x, skip], dim=1)
103
104         x = self.conv_block(x)
105         return x
106
107
108 class UNetDecoder(nn.Module):
109     """
110     The UNet decoder.
111
112     Args:
113         encoder_hidden_dim (int): The hidden dimension of the encoder.
114         decoder_channels (list): A list of the number of channels for each decoder block.
115
116     Returns:

```

```

117         x (torch.Tensor): The output tensor after the decoder.
118     """
119     def __init__(self, encoder_hidden_dim, decoder_channels):
120         super().__init__()
121
122         self.init_conv = nn.Conv2d(encoder_hidden_dim, decoder_channels[0], kernel_size=1)
123
124         self.decoder_blocks = nn.ModuleList()
125         in_channels = decoder_channels[0]
126         for i in range(len(decoder_channels)-1):
127             out_ch = decoder_channels[i+1]
128
129             block = DecoderBlock(
130                 in_channels=in_channels,
131                 in_channels_skip=encoder_hidden_dim,
132                 out_channels=out_ch
133             )
134             self.decoder_blocks.append(block)
135
136             in_channels = out_ch
137
138
139     def forward(self, x, skips):
140         x = self.init_conv(x) # 768 -> 1024 channels
141         for block, skip in zip(self.decoder_blocks, reversed(skips)):
142             x = block(x, skip)
143
144         return x
145
146
147     class ClipUNet(nn.Module):
148         """
149         The main ClipUNet model.
150
151         Args:
152             num_classes (int): The number of output classes.
153             decoder_channels (list): A list of the number of channels for each decoder block.
154             freeze_encoder (bool): Whether to freeze the CLIP encoder's parameters.
155             model_name (str): The name of the pre-trained CLIP model to use.
156             skip_indices (list): The indices of the hidden states to use for skip connections.
157
158         Returns:
159             output (torch.Tensor): The output tensor.
160         """
161         def __init__(self,
162             num_classes=4,
163             decoder_channels=[1024, 512, 256, 128, 64],
164             freeze_encoder=True,
165             model_name="openai/clip-vit-base-patch16",
166             skip_indices=[3, 5, 7, 9]
167         ):
168             super().__init__()
169
170             self.encoder = ClipViTEncoder(
171                 model_name=model_name,
172                 freeze_encoder=freeze_encoder,
173                 skip_indices=skip_indices
174             )
175
176             self.decoder = UNetDecoder(
177                 encoder_hidden_dim=self.encoder.hidden_dim,
178                 decoder_channels=decoder_channels
179             )
180
181             self.output_layer = nn.Conv2d(decoder_channels[-1], num_classes, kernel_size=1)
182
183

```



```

184     def forward(self, x):
185         x, skips = self.encoder(x)
186         decoder_output = self.decoder(x, skips)
187         output = self.output_layer(decoder_output)
188         return output
189

```

Clip without skip connections (clip/clipunet_{noskips}.py)

```

1  import torch.nn as nn
2  from transformers import CLIPVisionModel, CLIPVisionConfig
3
4  PRETRAINED_MODEL_NAME = "openai/clip-vit-base-patch16"
5
6  class ClipViTEncoderNoSkips(nn.Module):
7      """
8      This class implements a CLIP ViT encoder without skip connections.
9
10     Args:
11         model_name (str, optional): The name of the pre-trained CLIP model to use. Defaults to "openai/clip-vit-base-patch16".
12         freeze_encoder (bool, optional): Whether to freeze the encoder weights. Defaults to True.
13
14     Returns:
15         torch.Tensor: The bottleneck features of the input image.
16     """
17     def __init__(self, model_name="openai/clip-vit-base-patch16", freeze_encoder=True): # Removed bottleneck_index
18         super().__init__()
19
20         self.config = CLIPVisionConfig.from_pretrained(model_name)
21         self.clip_vit = CLIPVisionModel.from_pretrained(model_name)
22
23         if freeze_encoder:
24             # Freeze the encoder parameters if specified
25             for param in self.clip_vit.parameters():
26                 param.requires_grad = False
27
28         self.grid_size = self.config.image_size // self.config.patch_size
29         self.hidden_dim = self.config.hidden_size
30
31     def forward(self, x):
32         if x.shape[2] != self.config.image_size or x.shape[3] != self.config.image_size:
33             # Warn the user if the input image size doesn't match the expected size
34             print(
35                 f"Input image size ({x.shape[2]}x{x.shape[3]}) doesn't match "
36                 f"CLIP expected size ({self.config.image_size}x{self.config.image_size}). "
37                 f"Behavior may be unexpected. Consider resizing input."
38             )
39
40         outputs = self.clip_vit(pixel_values=x)
41         last_hidden_state = outputs.last_hidden_state
42
43         patch_embeddings = last_hidden_state[:, 1:, :] # Remove CLS token
44         bottleneck_features = patch_embeddings \
45             .reshape(x.shape[0], self.grid_size, self.grid_size, self.hidden_dim) \
46             .permute(0, 3, 1, 2).contiguous()
47
48         return bottleneck_features
49
50
51     class DecoderBlockNoSkip(nn.Module):
52         """
53         This class implements a decoder block without skip connections.
54
55     Args:
56         in_channels (int): Number of input channels.
57         out_channels (int): Number of output channels.

```

```

58
59     Returns:
60         torch.Tensor: Output tensor after the decoder block.
61     """
62     def __init__(self, in_channels, out_channels):
63         super().__init__()
64
65         # Upsample the input feature map
66         self.upsample = nn.ConvTranspose2d(in_channels, in_channels, kernel_size=2, stride=2) # Maybe reduce channels here
67
68         # Convolutional block
69         self.conv_block = nn.Sequential(
70             nn.Conv2d(in_channels, out_channels, kernel_size=3, padding=1, bias=False), # Adjusted input channels here
71             nn.BatchNorm2d(out_channels),
72             nn.ReLU(inplace=True),
73             nn.Conv2d(out_channels, out_channels, kernel_size=3, padding=1, bias=False),
74             nn.BatchNorm2d(out_channels),
75             nn.ReLU(inplace=True),
76         )
77
78     def forward(self, x):
79         x = self.upsample(x)
80         x = self.conv_block(x)
81         return x
82
83
84 class UNetDecoderNoSkips(nn.Module):
85     """
86     This class implements a UNet decoder without skip connections.
87
88     Args:
89         encoder_hidden_dim (int): The hidden dimension of the encoder.
90         decoder_channels (list of int): The number of channels for each decoder block.
91
92     Returns:
93         torch.Tensor: The output tensor after passing through the decoder.
94     """
95     def __init__(self, encoder_hidden_dim, decoder_channels):
96         super().__init__()
97
98         # Initial convolution layer
99         self.init_conv = nn.Conv2d(encoder_hidden_dim, decoder_channels[0], kernel_size=1)
100
101         # Decoder blocks
102         self.decoder_blocks = nn.ModuleList()
103         in_channels = decoder_channels[0]
104         for i in range(len(decoder_channels)-1):
105             out_ch = decoder_channels[i+1]
106
107             block = DecoderBlockNoSkip(
108                 in_channels=in_channels,
109                 out_channels=out_ch,
110             )
111             self.decoder_blocks.append(block)
112             in_channels = out_ch
113
114     def forward(self, x):
115         x = self.init_conv(x)
116         for block in self.decoder_blocks:
117             x = block(x)
118         return x
119
120
121 class ClipUNetNoSkips(nn.Module):
122     """
123     This class implements a CLIP UNet without skip connections.
124

```

```

125     Args:
126         num_classes (int, optional): The number of output classes. Defaults to 4.
127         decoder_channels (list of int, optional): The number of channels for each decoder block. Defaults to [1024, 512, 256, 128, 64].
128         freeze_encoder (bool, optional): Whether to freeze the encoder weights. Defaults to True.
129         model_name (str, optional): The name of the pre-trained CLIP model to use. Defaults to "openai/clip-vit-base-patch16".
130
131     Returns:
132         torch.Tensor: The output tensor after passing through the UNet.
133     """
134     def __init__(self,
135                 num_classes=4,
136                 decoder_channels=[1024, 512, 256, 128, 64],
137                 freeze_encoder=True,
138                 model_name="openai/clip-vit-base-patch16"
139                 ):
140         super().__init__()
141
142         # Encoder
143         self.encoder = ClipViTEncoderNoSkips(
144             model_name=model_name,
145             freeze_encoder=freeze_encoder
146         )
147
148         # Decoder
149         self.decoder = UNetDecoderNoSkips(
150             encoder_hidden_dim=self.encoder.hidden_dim,
151             decoder_channels=decoder_channels
152         )
153         # Output layer
154         self.output_layer = nn.Conv2d(decoder_channels[-1], num_classes, kernel_size=1)
155
156     def forward(self, x):
157         x = self.encoder(x)
158         decoder_output = self.decoder(x)
159         output = self.output_layer(decoder_output)
160         return output

```

Prompt based model (prompt_{based}/prompt.py)

```

1  from clip.clipunet import ClipUNet
2  from unet.unet import unet
3  import torch
4  from torch import nn
5
6  class PromptModel(nn.Module):
7      """
8      Initializes the PromptModel.
9      Args:
10         path (str, optional): Path to the checkpoint file. Defaults to None.
11     """
12     def __init__(self, path=None):
13         super().__init__()
14
15         self.clip = ClipUNet()
16         self.mask = unet(4, 1)
17         self.softmax = nn.Softmax(dim=1)
18         self.sigmoid = nn.Sigmoid()
19
20         if path is not None:
21             try:
22                 # Load the checkpoint.
23                 checkpoint = torch.load(path, weights_only=False, map_location=lambda storage, loc: storage) # Load to CPU
24                 self.clip.load_state_dict(checkpoint["model_state_dict"])
25             except Exception as e:
26                 print(f"Error loading checkpoint: {str(e)[:200]}")
27                 raise

```

```

28
29     # Freeze the clip parameters.
30     for param in self.clip.parameters():
31         param.requires_grad = False
32
33     def forward(self, x, heatmap):
34         # Pass input through clip model.
35         clip_logit = self.clip(x)
36         # Apply softmax to clip logits.
37         clip_prob = self.softmax(clip_logit)
38
39         # Concatenate input and heatmap, pass through mask model.
40         mask_logit = self.mask(torch.concat([x, heatmap], dim=1))
41         # Apply sigmoid to mask logits.
42         mask_prob = self.sigmoid(mask_logit)
43
44         # Initialize tensor to store final probabilities.
45         final_probs = torch.empty_like(clip_prob)
46         # Calculate selected probabilities.
47         selected_prob = mask_prob * clip_prob
48
49         # Assign probabilities to the final tensor.
50         # Assign the selected probabilities for background, cat, and dog.
51         final_probs[:, 1:4, :, :] = selected_prob[:, 0:3, :, :]
52         # Assign the mask probability for deactivated class.
53         final_probs[:, 0:1, :, :] = 1.0 - mask_prob
54         # Merge boundary class with background
55         final_probs[:, 1:2, :, :] += selected_prob[:, 3:4, :, :]
56
57     return final_probs

```

Datasets (utils/dataset.py)

```

1  import os
2  import matplotlib.pyplot as plt
3  from torch.utils.data import Dataset
4  from torchvision.io import decode_image
5
6  class dataset(Dataset):
7      """
8      Initializes the dataset class.
9      Args:
10         img_dir (str): Directory containing the images.
11         label_dir (str): Directory containing the labels.
12         transform (callable, optional): Optional transform to be applied on a sample. Defaults to None.
13         target_transform (callable, optional): Optional transform to be applied on a target. Defaults to None.
14      """
15     def __init__(self, img_dir, label_dir, transform=None, target_transform=None):
16         # Initialize the image and label directories
17         self.img_dir = img_dir
18         self.label_dir = label_dir
19         # Get the names of the images, without the extension, and sort them
20         self.img_names = sorted([os.path.splitext(filename)[0] for filename in os.listdir(img_dir)])
21         # Calculate the length of the dataset
22         self.len = len(self.img_names)
23         # Initialize the transforms
24         self.transform = transform
25         self.target_transform = target_transform
26
27     def __len__(self):
28         return self.len
29
30     def __getitem__(self, idx):
31         """
32         Loads and returns an image label pair from the dataset.
33         Args:

```



```

34         idx (int): Index of the item to retrieve.
35     Returns:
36         tuple: A tuple containing the image and its corresponding label.
37     """
38     # Load the image and normalize it
39     img = decode_image(os.path.join(self.img_dir, self.img_names[idx] + ".jpg")).float() / 255
40     # Load the label
41     label = decode_image(os.path.join(self.label_dir, self.img_names[idx] + ".png"))
42
43     # Apply transformations to the image if specified
44     if self.transform:
45         img = self.transform(img)
46
47     # Apply transformations to the label if specified
48     if self.target_transform:
49         label = self.target_transform(label)
50
51     return img, label
52
53 class promptDataset(Dataset):
54     """
55     Initializes the promptDataset class.
56     Args:
57         img_dir (str): Directory containing the images.
58         heatmap_dir (str): Directory containing the heatmaps.
59         label_dir (str): Directory containing the labels.
60         transform (callable, optional): Optional transform to be applied on a sample. Defaults to None.
61         target_transform (callable, optional): Optional transform to be applied on a target. Defaults to None.
62     """
63     def __init__(self, img_dir, heatmap_dir, label_dir, transform=None, target_transform=None):
64         # Initialize the image, heatmap, and label directories
65         self.img_dir = img_dir
66         self.heatmap_dir = heatmap_dir
67         self.label_dir = label_dir
68         # Get the names of the images, without the extension, and sort them
69         self.img_names = sorted([os.path.splitext(filename)[0] for filename in os.listdir(img_dir)])
70         # Calculate the length of the dataset
71         self.len = len(self.img_names)
72         # Initialize the transforms
73         self.transform = transform
74         self.target_transform = target_transform
75
76     def __len__(self):
77         return self.len
78
79     def __getitem__(self, idx):
80         """
81         Loads and returns an image, heatmap, label triplet from the dataset.
82         Args:
83             idx (int): Index of the item to retrieve.
84         Returns:
85             tuple: A tuple containing the image, heatmap, and its corresponding label.
86         """
87         # Load the image and normalize it
88         img = decode_image(os.path.join(self.img_dir, self.img_names[idx] + ".jpg")).float()/255
89         # Load the heatmap and normalize it
90         heatmap = decode_image(os.path.join(self.heatmap_dir, self.img_names[idx] + ".png"))/255
91         # Load the label
92         label = decode_image(os.path.join(self.label_dir, self.img_names[idx] + ".png"))
93
94         # Apply transformations to the image if specified
95         if self.transform:
96             img = self.transform(img)
97
98         # Apply transformations to the label if specified
99         if self.target_transform:
100             label = self.target_transform(label)

```

```

101
102
103         return img, heatmap, label
104
105
106 def display_img_label(data, idx):
107     """
108     Displays an image and its corresponding label using matplotlib.
109     Args:
110         data (torch.utils.data.Dataset): The dataset containing images and labels.
111         idx (int): The index of the image and label to display.
112     """
113     # Get the image and label from the dataset
114     img, label = data[idx]
115     # Create a figure with two subplots
116     figure = plt.figure(figsize=(10, 20))
117     # Add the first subplot for the image
118     figure.add_subplot(1, 2, 1)
119     # Display the image, permuting the dimensions to match matplotlib's expected format
120     plt.imshow(img.permute(1, 2, 0))
121
122     # Add the second subplot for the label
123     figure.add_subplot(1, 2, 2)
124     # Display the label as a grayscale image, permuting the dimensions
125     plt.imshow(label.permute(1, 2, 0), cmap='grey')
126
127     # Show the plot
128     plt.show()
129
130
131 class target_remap(object):
132     """
133     Remaps boundary class (255) to 3
134     """
135     def __call__(self, img):
136         # Remap pixel value 255 to 3
137         img[img == 255] = 3
138         return img
139
140
141 def diff_size_collate(batch):
142     """
143     Collates a batch of data with potentially different image sizes.
144     Args:
145         batch (list): A list of tuples, where each tuple contains an image and its label.
146     Returns:
147         tuple: A tuple containing two lists: a list of images and a list of labels.
148     """
149     # Extract images from the batch
150     imgs = [item[0] for item in batch]
151     # Extract labels from the batch
152     labels = [item[1] for item in batch]
153     return imgs, labels

```

MetricsHistory (utils/MetricsHistory.py)

```

1 import torch
2 import torch.nn.functional as F
3
4 class MetricsHistory:
5     """
6     Accumulates TP, FP, FN, TN over an epoch for multi-class segmentation
7     and computes Dice, IoU, and Accuracy metrics.
8     """
9     def __init__(self, num_classes: int, ignore_index: int = None, device: str = 'cpu'):
10         """

```

```

11     Initializes the MetricsHistory object.
12     Args:
13         num_classes (int): Number of classes including background.
14         ignore_index (int, optional): Index of the class to ignore during metric calculation. Defaults to None.
15         device (str): Device to perform initial calculations, results accumulated on CPU.
16     """
17     self.num_classes = num_classes
18     self.ignore_index = ignore_index
19
20     # Initialize tensors to store the sums of TP, FP, FN, and TN.
21     self.total_tp = torch.zeros(num_classes, dtype=torch.float64, device='cpu')
22     self.total_fp = torch.zeros(num_classes, dtype=torch.float64, device='cpu')
23     self.total_fn = torch.zeros(num_classes, dtype=torch.float64, device='cpu')
24     self.total_tn = torch.zeros(num_classes, dtype=torch.float64, device='cpu')
25
26     # History lists for epoch metrics
27     self.epoch_mean_dice_history = []
28     self.epoch_mean_iou_history = []
29     self.epoch_mean_acc_history = []
30
31     self.epoch_per_class_dice_history = []
32     self.epoch_per_class_iou_history = []
33     self.epoch_per_class_acc_history = []
34
35     self.last_per_class_iou = None
36     self.last_per_class_dice = None
37     self.last_per_class_acc = None
38
39     # Metric mask (calculated once)
40     self.mask = torch.ones(num_classes, dtype=torch.bool)
41     # Set mask element to False if ignore_index is specified.
42     if self.ignore_index is not None and 0 <= self.ignore_index < self.num_classes:
43         self.mask[self.ignore_index] = False
44
45 def reset(self):
46     """
47     Resets the accumulated TP, FP, FN, TN counts.
48     """
49     # Resets the accumulated statistics to zero at the start of each epoch.
50     self.total_tp.zero_()
51     self.total_fp.zero_()
52     self.total_fn.zero_()
53     self.total_tn.zero_()
54
55 def accumulate(self, pred: torch.Tensor, label: torch.Tensor):
56     """
57     Accumulates statistics for a single prediction-label pair.
58
59     Args:
60         pred (torch.Tensor): Predicted logits or probabilities (C, H, W). Should be on self.device or moved.
61         label (torch.Tensor): Ground truth label map (H, W), LongTensor. Should be on self.device or moved.
62     """
63
64     # Get hard predictions
65     pred_hard = torch.argmax(pred.squeeze(0), dim=0) # (H, W)
66
67     # One-hot encode
68     label_onehot = F.one_hot(label.squeeze(0), num_classes=self.num_classes).permute(2, 0, 1).bool() # (C, H, W)
69     pred_onehot = F.one_hot(pred_hard, num_classes=self.num_classes).permute(2, 0, 1).bool() # (C, H, W)
70
71     # Calculate TP, FP, FN, TN per class
72     tp = (pred_onehot & label_onehot).sum(dim=(1, 2))
73     fp = (pred_onehot & ~label_onehot).sum(dim=(1, 2))
74     fn = (~pred_onehot & label_onehot).sum(dim=(1, 2))
75     tn = (~pred_onehot & ~label_onehot).sum(dim=(1, 2))
76
77     # tp = (pred_onehot & label_onehot).sum(dim=(1, 2))

```

```

78     # fp = pred_onehot.sum(dim=(1, 2)) - tp
79     # fn = label_onehot.sum(dim=(1, 2)) - tp
80     # tn = label.numel() - fn - fp - tp
81
82     # Accumulate on CPU with float64
83     self.total_tp += tp.cpu().to(torch.float64)
84     self.total_fp += fp.cpu().to(torch.float64)
85     self.total_fn += fn.cpu().to(torch.float64)
86     self.total_tn += tn.cpu().to(torch.float64) # Accumulate TN if needed for accuracy
87
88
89     def compute_epoch_metrics(self, epsilon: float = 1e-6):
90         """
91         Computes the macro-averaged metrics for the accumulated epoch statistics,
92         appends them to the history lists, and returns the computed mean metrics.
93
94         Args:
95             epsilon (float): Small value to avoid division by zero.
96
97         Returns:
98             tuple: (mean_dice, mean_iou, mean_acc) for the current epoch.
99         """
100
101         tp = self.total_tp
102         fp = self.total_fp
103         fn = self.total_fn
104         tn = self.total_tn
105
106         per_class_iou = tp / (tp + fp + fn)
107         per_class_dice = (2 * tp) / (2 * tp + fp + fn)
108         per_class_acc = (tp + tn) / (tp + tn + fp + fn)
109
110         # Compute the mean IoU, Dice, and accuracy, considering the mask.
111         mean_iou = per_class_iou[self.mask].mean().item()
112         mean_dice = per_class_dice[self.mask].mean().item()
113         mean_acc = per_class_acc[self.mask].mean().item()
114
115         # Append to history
116         self.epoch_mean_iou_history.append(mean_iou)
117         self.epoch_mean_dice_history.append(mean_dice)
118         self.epoch_mean_acc_history.append(mean_acc)
119
120         self.epoch_per_class_iou_history.append(per_class_iou.numpy())
121         self.epoch_per_class_dice_history.append(per_class_dice.numpy())
122         self.epoch_per_class_acc_history.append(per_class_acc.numpy())
123
124         self.last_per_class_iou = per_class_iou
125         self.last_per_class_dice = per_class_dice
126         self.last_per_class_acc = per_class_acc
127
128         return mean_dice, mean_iou, mean_acc
129
130     def to(self, device):
131         """
132         Moves the internal tensors to the specified device.
133
134         Args:
135             device (str): The device to move the tensors to (e.g., 'cuda', 'cpu').
136         """
137         self.total_tp = self.total_tp.to(device)
138         self.total_fp = self.total_fp.to(device)
139         self.total_fn = self.total_fn.to(device)
140         self.total_tn = self.total_tn.to(device)
141
142         self.mask = self.mask.to(device)
143
144         if self.last_per_class_iou is not None:
145             self.last_per_class_iou = self.last_per_class_iou.to(device)

```

```

145         if self.last_per_class_dice is not None:
146             self.last_per_class_dice = self.last_per_class_dice.to(device)
147
148         if self.last_per_class_acc is not None:
149             self.last_per_class_acc = self.last_per_class_acc.to(device)
150
151     def get_ignore_index(self):
152         return self.ignore_index
153
154     def get_num_classes(self):
155         return self.num_classes
156
157     def get_mean_dice_history(self):
158         return self.epoch_mean_dice_history
159
160     def get_mean_iou_history(self):
161         return self.epoch_mean_iou_history
162
163     def get_mean_acc_history(self):
164         return self.epoch_mean_acc_history
165
166     def get_class_dice_history(self):
167         return self.epoch_per_class_dice_history
168
169     def get_class_iou_history(self):
170         return self.epoch_per_class_iou_history
171
172     def get_class_acc_history(self):
173         return self.epoch_per_class_acc_history
174
175     def get_last_per_class_dice(self):
176         return self.last_per_class_dice
177
178     def get_last_per_class_iou(self):
179         return self.last_per_class_iou
180
181     def get_last_per_class_acc(self):
182         return self.last_per_class_acc
183

```

Training Pipeline (utils/training.py)

```

1  import os
2  import numpy as np
3  import torch
4  from torch import nn
5  from torch.utils.data import DataLoader
6  from tqdm.notebook import tqdm
7  from utils.utils import process_batch_forward, process_batch_reverse
8  from utils.MetricsHistory import MetricsHistory
9  from torchvision.transforms import InterpolationMode
10
11  if torch.backends.mps.is_available():
12      device = torch.device("mps")
13  elif torch.cuda.is_available():
14      device = torch.device("cuda")
15  else:
16      device = torch.device("cpu")
17
18  def train_loop(dataloader, model, loss_fn, optimizer, accumulation_steps, device, scheduler=None, target_size = None):
19      """
20      Performs one epoch of training.
21      Args:
22          dataloader: DataLoader for training data.
23          model: The neural network model.
24          loss_fn: The loss function.

```

```

25     optimizer: The optimizer.
26     accumulation_steps: Number of steps to accumulate gradients before updating.
27     device: The device to train on (CPU or GPU).
28     scheduler: Learning rate scheduler (optional).
29     target_size: Target size for image resizing (optional).
30     """
31     model.train()
32     total_loss = 0.0
33     processed_batches = 0
34
35     optimizer.zero_grad()
36
37     pbar = tqdm(enumerate(dataloader), total=len(dataloader), desc="Training")
38     for batch_idx, (X, y) in pbar:
39
40         if target_size is not None:
41             # Resize images and labels if target_size is specified
42             X, _ = process_batch_forward(X, target_size=target_size)
43             y, _ = process_batch_forward(y, target_size=target_size, interpolation=InterpolationMode.NEAREST)
44
45             X, y = X.to(device), y.to(device).long()
46             pred = model(X)
47             loss = loss_fn(pred, y.squeeze(1))
48
49             scaled_loss = loss / accumulation_steps
50             scaled_loss.backward()
51
52             if (batch_idx + 1) % accumulation_steps == 0 or (batch_idx + 1) == len(dataloader):
53                 optimizer.step()
54                 if scheduler:
55                     scheduler.step()
56                 optimizer.zero_grad()
57
58                 total_loss += loss.item()
59                 processed_batches += 1
60                 pbar.set_postfix({'loss': loss.item(), 'lr': optimizer.param_groups[0]['lr']})
61
62     avg_loss = total_loss / processed_batches if processed_batches > 0 else 0
63     print(f"Training Avg loss (per effective batch): {avg_loss:>8f}")
64     return avg_loss
65
66
67 def eval_loop(dataloader, model, loss_fn, device, target_size, agg):
68     """
69     Evaluation loop calculating loss, aggregated Dice, aggregated Acc, and aggregated IoU.
70
71     Args:
72         dataloader: yields batches of (list[Tensor(C,H,W)], list[Tensor(H,W)])
73         model: the neural network model (on device)
74         loss_fn: the combined loss function (e.g., DiceCELoss) used for training
75         device: the torch device (cuda or cpu)
76         target_size: the size the model expects for input
77     """
78     model.eval()
79     num_images_processed = 0
80     total_loss = 0.0
81     num_classes = 4
82     agg.reset()
83
84     with torch.no_grad():
85         for X, y in tqdm(dataloader, desc="Eval"):
86
87             X, meta_list = process_batch_forward(X, target_size=target_size)
88             X = X.to(device)
89             preds = model(X) # Logits [N, C, target, target]
90
91             preds = process_batch_reverse(preds, meta_list, interpolation='bilinear')

```



```

92
93     for pred, label in zip(preds, y):
94         pred = pred.to(device) # (C,H,W)
95         label = label.to(device).long() # (H,W)
96
97         loss = loss_fn(pred.unsqueeze(0), label.unsqueeze(0).squeeze(1)) # Add batch dimension
98         total_loss += loss.item()
99         agg.accumulate(pred, label)
100
101         num_images_processed += 1
102
103     avg_loss = total_loss / num_images_processed
104
105     mean_dice, mean_iou, mean_acc = agg.compute_epoch_metrics()
106     per_class_iou = agg.get_last_per_class_iou()
107     ignore_index = agg.get_ignore_index()
108
109     print(f"\n--- Evaluation Complete ---")
110     print(f" Images Processed: {num_images_processed}")
111     print(f" Average Loss (Original Size): {avg_loss:>8f}")
112     print(f" Ignored Class : {ignore_index}")
113     print(f" Macro Avg Acc score: {mean_acc:>8f}")
114     print(f" Macro Avg Dice Score: {mean_dice:>8f}")
115     print(f" Mean IoU (mIoU): {mean_iou:>8f}")
116     print(f" --- Per-Class IoU ---")
117     for c in range(num_classes):
118         print(f" Class {c}: {per_class_iou[c].item():>8f}")
119     print("-" * 25)
120
121     return avg_loss, mean_dice, mean_iou
122
123 def trainReconstruction(dataloader, model, loss_fn, optimizer, accumulation_steps):
124     """
125     Trains a reconstruction model.
126     Args:
127         dataloader: DataLoader for the training data.
128         model: The reconstruction model.
129         loss_fn: The loss function.
130         optimizer: The optimizer.
131         accumulation_steps: Number of steps to accumulate gradients.
132     """
133     losses = []
134     model.train()
135     for batch_idx, (X, _) in enumerate(tqdm(dataloader, total=len(dataloader), desc="Training")):
136         X = X.to(device)
137         # Compute prediction
138         pred = model(X)
139
140         # Compute loss
141         loss = loss_fn(pred, X)
142         losses.append(loss.item())
143         scaled_loss = loss / accumulation_steps
144
145         scaled_loss.backward()
146
147         if (batch_idx + 1) % accumulation_steps == 0 or (batch_idx + 1) == len(dataloader):
148             optimizer.step()
149             optimizer.zero_grad()
150
151     return np.mean(losses)
152
153 def train_loop_prompt(dataloader, model, loss_fn, optimizer, accumulation_steps, device, scheduler=None, target_size =
154     """
155     Performs one epoch of training for a prompt-based model.
156     Args:
157         dataloader: DataLoader for training data.
158         model: The prompt-based model.

```

```

159     loss_fn: The loss function.
160     optimizer: The optimizer.
161     accumulation_steps: Number of steps to accumulate gradients.
162     device: The device to train on (CPU or GPU).
163     scheduler: Learning rate scheduler (optional).
164     target_size: Target size for image resizing (optional).
165     """
166     model.train()
167     total_loss = 0.0
168     processed_batches = 0
169
170     optimizer.zero_grad()
171
172     pbar = tqdm(enumerate(dataloader), total=len(dataloader), desc="Training")
173     for batch_idx, (X, p, y) in pbar:
174
175         if target_size is not None:
176             X, _ = process_batch_forward(X, target_size=target_size)
177             p, _ = process_batch_forward(p, target_size=target_size)
178             y, _ = process_batch_forward(y, target_size=target_size, interpolation=InterpolationMode.NEAREST)
179
180             X, p, y = X.to(device), p.to(device), y.to(device).long()
181             pred = model(X, p)
182             loss = loss_fn(pred, y.squeeze(1))
183
184             scaled_loss = loss / accumulation_steps
185             scaled_loss.backward()
186
187             if (batch_idx + 1) % accumulation_steps == 0 or (batch_idx + 1) == len(dataloader):
188                 optimizer.step()
189                 if scheduler:
190                     scheduler.step()
191                 optimizer.zero_grad()
192
193                 total_loss += loss.item()
194                 processed_batches += 1
195                 pbar.set_postfix({'loss': loss.item(), 'lr': optimizer.param_groups[0]['lr']})
196
197     avg_loss = total_loss / processed_batches if processed_batches > 0 else 0
198     print(f"Training Avg loss (per effective batch): {avg_loss:>8f}")
199     return avg_loss
200
201
202 def evalReconstruction(dataloader, model, loss_fn, target_size, interpolation = 'bilinear'):
203     """
204     Evaluates a reconstruction model.
205     Args:
206         dataloader: DataLoader for the evaluation data.
207         model: The reconstruction model.
208         loss_fn: The loss function.
209         target_size: Target size for image resizing.
210         interpolation: Interpolation mode for resizing.
211     """
212     model.eval()
213     num_batches = len(dataloader)
214     total_loss = 0.0
215     losses = []
216     with torch.no_grad():
217         for batch, (original_X, _) in enumerate(tqdm(dataloader, total=len(dataloader), desc="Evaluation")):
218             resized_X, meta_list = process_batch_forward(original_X, target_size=target_size) # resize X for network
219             resized_X = resized_X.to(device)
220
221             # Compute prediction
222             pred = model(resized_X)
223
224             pred = process_batch_reverse(pred, meta_list, interpolation=interpolation)
225

```

```

226         for p, label in zip(pred, original_X):
227             # Move individual prediction and label to the device
228             p = p.to(device).unsqueeze(0) # Add batch dimension
229             label = label.to(device).unsqueeze(0) # Add batch dimension and ensure type is long
230
231             if label.shape[1] == 4 and label.ndim == 4:
232                 label = label[:, :3, :, :] # RGBA to RGB
233
234             loss = loss_fn(p, label.squeeze(1))
235             total_loss += loss.item()
236             # Loss list
237             losses.append(loss.item())
238
239         return total_loss / num_batches, np.mean(losses)
240
241
242 def eval_loop_prompt(dataloader, model, loss_fn, device, target_size, agg):
243     """
244     Evaluation loop calculating loss, aggregated Dice, and aggregated IoU for a prompt-based model.
245
246     Args:
247         dataloader: yields batches of (list[Tensor(C,H,W)], list[Tensor(H,W)])
248         model: the neural network model (on device)
249         loss_fn: the combined loss function (e.g., DiceCELoss) used for training
250         device: the torch device (cuda or cpu)
251         target_size: the size the model expects for input
252     """
253     model.eval()
254     num_images_processed = 0
255     total_loss = 0.0
256     num_classes = 4
257
258     with torch.no_grad():
259         for X, p, y in tqdm(dataloader, desc="Eval"):
260
261             X, meta_list = process_batch_forward(X, target_size=target_size)
262             p, _ = process_batch_forward(p, target_size=target_size)
263             X, p = X.to(device), p.to(device)
264             preds = model(X, p) # Logits [N, C, target, target]
265
266             preds = process_batch_reverse(preds, meta_list, interpolation='bilinear')
267
268             for pred, label in zip(preds, y):
269                 pred = pred.to(device) # (C,H,W)
270                 label = label.to(device).long() # (H,W)
271
272                 loss = loss_fn(pred.unsqueeze(0), label.unsqueeze(0).squeeze(1)) # Add batch dimension
273                 total_loss += loss.item()
274                 agg.accumulate(pred, label)
275
276                 num_images_processed += 1
277
278     avg_loss = total_loss / num_images_processed
279
280     mean_dice, mean_iou, mean_acc = agg.compute_epoch_metrics()
281     per_class_iou = agg.get_last_per_class_iou()
282     ignore_index = agg.get_ignore_index()
283
284     print(f"\n--- Evaluation Complete ---")
285     print(f" Images Processed: {num_images_processed}")
286     print(f" Average Loss (Original Size): {avg_loss:>8f}")
287     print(f" Ignored Class : {ignore_index}")
288     print(f" Macro Avg Acc score: {mean_acc:>8f}")
289     print(f" Macro Avg Dice Score: {mean_dice:>8f}")
290     print(f" Mean IoU (mIoU): {mean_iou:>8f}")
291     print(f" --- Per-Class IoU ---")
292     for c in range(num_classes):

```

```

293         print(f"      Class {c}: {per_class_iou[c].item():>8f}")
294     print("-" * 25)
295
296     return avg_loss, mean_dice, mean_iou
297
298
299 def start_prompt(
300     model_save_dir: str,
301     model_save_name: str,
302     model: nn.Module,
303     optimizer: torch.optim,
304     train_dataloader: DataLoader,
305     val_dataloader: DataLoader,
306     accumulation_steps: int,
307     device: torch.device,
308     train_loss_fn: nn.Module,
309     val_loss_fn: nn.Module,
310     target_size: int,
311     scheduler: torch.optim.lr_scheduler = None,
312     agg: MetricsHistory = None,
313     load: bool = True,
314     save: bool = True,
315     num_classes: int = 4,
316     ignore_index: int = 3,
317     epochs: int = 100,
318 ):
319     """
320     Starts the training pipeline for a prompt-based model.
321     Args:
322         model_save_dir: Directory to save the model.
323         model_save_name: Name of the model file.
324         model: The prompt-based model.
325         optimizer: The optimizer.
326         train_dataloader: DataLoader for the training data.
327         val_dataloader: DataLoader for the validation data.
328         accumulation_steps: Number of steps to accumulate gradients.
329         device: The device to train on (CPU or GPU).
330         train_loss_fn: The loss function for training.
331         val_loss_fn: The loss function for validation.
332         target_size: Target size for image resizing.
333         scheduler: Learning rate scheduler (optional).
334         agg: MetricsHistory object (optional).
335         load: Whether to load a checkpoint (default: True).
336         save: Whether to save the model (default: True).
337         num_classes: Number of classes (default: 4).
338         ignore_index: Index to ignore in the loss calculation (default: 3).
339         epochs: Number of training epochs (default: 100).
340     """
341     start_epoch = 0
342     best_dev_dice = -np.inf
343     best_dev_miou = -np.inf
344     best_dev_loss = -np.inf
345
346     os.makedirs(model_save_dir, exist_ok=True)
347     os.makedirs(f"{model_save_dir}/metrics", exist_ok=True)
348     if load and os.path.isfile(f"{model_save_dir}/{model_save_name}"):
349         print(f"Loading checkpoint from: {model_save_dir}/{model_save_name}")
350
351         checkpoint = torch.load(f"{model_save_dir}/{model_save_name}", map_location=device)
352
353         model.load_state_dict(checkpoint["model_state_dict"])
354         print(" -> Model state loaded.")
355
356         # Load optimizer state
357         try:
358             optimizer.load_state_dict(checkpoint["optimizer_state_dict"])
359             print(" -> Optimizer state loaded.")

```

```

360     except Exception as e:
361         print(f" -> Warning: Could not load optimizer state: {e}. Optimizer will start from scratch.")
362
363     # Load scheduler state
364     try:
365         scheduler.load_state_dict(checkpoint["scheduler_state_dict"])
366         print(" -> Scheduler state loaded.")
367     except Exception as e:
368         print(f" -> Warning: Could not load scheduler state: {e}. Scheduler will start from scratch.")
369
370     # Load Metrics History
371     try:
372         agg = checkpoint.get("history")
373         agg.to(device)
374         print(" -> Metrics History loaded.")
375     except Exception as e:
376         print(f" -> No metric history saved")
377         agg = MetricsHistory(num_classes, ignore_index)
378
379     # Load training metadata
380     start_epoch = checkpoint.get("epoch", 0)
381     best_dev_dice = checkpoint.get("best_dev_dice", -np.inf)
382     best_dev_miou = checkpoint.get("best_dev_miou", -np.inf)
383     best_dev_loss = checkpoint.get("best_dev_loss", np.inf)
384
385     print(f" -> Resuming training from epoch {start_epoch + 1}")
386     print(f" -> Loaded best metrics: Dice={best_dev_dice:.6f}, mIoU={best_dev_miou:.6f}, Loss={best_dev_loss:.6f}")
387     loaded_notes = checkpoint.get("notes", "N/A")
388     print(f" -> Notes from checkpoint: {loaded_notes}")
389
390 else:
391     print(f"Checkpoint file not found at {model_save_dir}/{model_save_name}. Starting training from scratch.")
392
393 # --- Training and Evaluation Loop ---
394 print("\nStarting Training...")
395 for t in range(start_epoch, epochs):
396     print(f"Epoch {t+1}\n-----")
397
398     train_loss = train_loop_prompt(train_dataloader, model, train_loss_fn, optimizer, accumulation_steps, device,
399     val_loss, val_dice, val_miou = eval_loop_prompt(val_dataloader, model, val_loss_fn, device, target_size, agg)
400
401     if save:
402         metrics = {
403             "epoch": t + 1,
404             "history": agg
405         }
406         torch.save(metrics, f"{model_save_dir}/metrics/{model_save_name}")
407
408     if val_miou > best_dev_miou:
409         best_dev_dice = val_dice
410         best_dev_miou = val_miou
411         best_dev_loss = val_loss
412
413     if save and scheduler:
414         print(f"Validation IoU score improved ({best_dev_miou:.6f}). Saving model...")
415
416         checkpoint = {
417             "epoch": t + 1,
418             "model_state_dict": model.state_dict(),
419             "optimizer_state_dict": optimizer.state_dict(),
420             "scheduler_state_dict": scheduler.state_dict(),
421             "best_dev_dice": best_dev_dice,
422             "best_dev_miou": best_dev_miou,
423             "best_dev_loss": best_dev_loss,
424             "history": agg,
425             "notes": f"Model saved based on best Micro Dice. Ignored index for metric: {ignore_index}"
426         }

```

```

427         torch.save(checkpoint, f"{model_save_dir}/{model_save_name}")
428     elif save:
429         print(f"Validation IoU score improved ({best_dev_miou:.6f}). Saving model...")
430
431         checkpoint = {
432             "epoch": t + 1,
433             "model_state_dict": model.state_dict(),
434             "optimizer_state_dict": optimizer.state_dict(),
435             "best_dev_dice": best_dev_dice,
436             "best_dev_miou": best_dev_miou,
437             "best_dev_loss": best_dev_loss,
438             "history": agg,
439             "notes": f"Model saved based on best Micro Dice. Ignored index for metric: {ignore_index}"
440         }
441         torch.save(checkpoint, f"{model_save_dir}/{model_save_name}")
442     else:
443         print(f"Validation IoU score did not improve from {best_dev_miou:.6f}")
444
445
446     print("\n--- Training Finished! ---")
447     print(f"Best validation IoU score achieved: {best_dev_miou:.6f}")
448     print(f"Corresponding validation dice: {best_dev_dice:.6f}")
449     print(f"Corresponding validation loss: {best_dev_loss:.6f}")
450     print(f"Best model saved to: {os.path.join(model_save_dir, model_save_name)}")
451
452
453 def start(
454     model_save_dir: str,
455     model_save_name: str,
456     model: nn.Module,
457     optimizer: torch.optim,
458     train_dataloader: DataLoader,
459     val_dataloader: DataLoader,
460     accumulation_steps: int,
461     device: torch.device,
462     train_loss_fn: nn.Module,
463     val_loss_fn: nn.Module,
464     target_size: int,
465     scheduler: torch.optim.lr_scheduler = None,
466     agg: MetricsHistory = None,
467     load: bool = True,
468     save: bool = True,
469     num_classes: int = 4,
470     ignore_index: int = 3,
471     epochs: int = 100,
472 ):
473     """
474     Starts the training pipeline.
475     Args:
476         model_save_dir: Directory to save the model.
477         model_save_name: Name of the model file.
478         model: The model.
479         optimizer: The optimizer.
480         train_dataloader: DataLoader for the training data.
481         val_dataloader: DataLoader for the validation data.
482         accumulation_steps: Number of steps to accumulate gradients.
483         device: The device to train on (CPU or GPU).
484         train_loss_fn: The loss function for training.
485         val_loss_fn: The loss function for validation.
486         target_size: Target size for image resizing.
487         scheduler: Learning rate scheduler (optional).
488         agg: MetricsHistory object (optional).
489         load: Whether to load a checkpoint (default: True).
490         save: Whether to save the model (default: True).
491         num_classes: Number of classes (default: 4).
492         ignore_index: Index to ignore in the loss calculation (default: 3).
493         epochs: Number of training epochs (default: 100).

```



```

494     """
495     start_epoch = 0
496     best_dev_dice = -np.inf
497     best_dev_miou = -np.inf
498     best_dev_loss = np.inf
499
500     os.makedirs(model_save_dir, exist_ok=True)
501     os.makedirs(f"{model_save_dir}/metrics", exist_ok=True)
502     if load and os.path.isfile(f"{model_save_dir}/{model_save_name}"):
503         print(f"Loading checkpoint from: {model_save_dir}/{model_save_name}")
504
505         # Load checkpoint
506         checkpoint = torch.load(f"{model_save_dir}/{model_save_name}", map_location=device, weights_only=True)
507
508         # Load model state
509         model.load_state_dict(checkpoint["model_state_dict"])
510         print(" -> Model state loaded.")
511
512         # Load optimizer state
513         try:
514             optimizer.load_state_dict(checkpoint["optimizer_state_dict"])
515             print(" -> Optimizer state loaded.")
516         except Exception as e:
517             print(f" -> Warning: Could not load optimizer state: {e}. Optimizer will start from scratch.")
518
519         # Load scheduler state
520         try:
521             scheduler.load_state_dict(checkpoint["scheduler_state_dict"])
522             print(" -> Scheduler state loaded.")
523         except Exception as e:
524             print(f" -> Warning: Could not load scheduler state: {e}. Scheduler will start from scratch.")
525
526         # Load Metrics History
527         try:
528             agg = checkpoint.get("history")
529             agg.to(device)
530             print(" -> Metrics History loaded.")
531         except Exception as e:
532             print(f" -> No metric history saved")
533             agg = MetricsHistory(num_classes, ignore_index)
534
535         # Load training metadata
536         start_epoch = checkpoint.get("epoch", 0)
537         best_dev_dice = checkpoint.get("best_dev_dice", -np.inf)
538         best_dev_miou = checkpoint.get("best_dev_miou", -np.inf)
539         best_dev_loss = checkpoint.get("best_dev_loss", np.inf)
540
541         print(f" -> Resuming training from epoch {start_epoch + 1}")
542         print(f" -> Loaded best metrics: Dice={best_dev_dice:.6f}, mIoU={best_dev_miou:.6f}, Loss={best_dev_loss:.6f}")
543         loaded_notes = checkpoint.get("notes", "N/A")
544         print(f" -> Notes from checkpoint: {loaded_notes}")
545
546     else:
547         print(f"Checkpoint file not found at {model_save_dir}/{model_save_name}. Starting training from scratch.")
548
549     # --- Training and Evaluation Loop ---
550     print("\nStarting Training...")
551     for t in range(start_epoch, epochs):
552         print(f"Epoch {t+1}\n-----")
553
554         train_loss = train_loop(train_dataloader, model, train_loss_fn, optimizer, accumulation_steps, device, scheduler)
555         val_loss, val_dice, val_miou = eval_loop(val_dataloader, model, val_loss_fn, device, target_size, agg)
556
557         if save:
558             metrics = {
559                 "epoch": t + 1,
560                 "history": agg

```

```

561         }
562         torch.save(metrics, f"{model_save_dir}/metrics/{model_save_name}")
563
564     if val_miou > best_dev_miou:
565         best_dev_dice = val_dice
566         best_dev_miou = val_miou
567         best_dev_loss = val_loss
568
569     if save and scheduler:
570         print(f"Validation IoU score improved ({best_dev_miou:.6f}). Saving model...")
571
572         checkpoint = {
573             "epoch": t + 1,
574             "model_state_dict": model.state_dict(),
575             "optimizer_state_dict": optimizer.state_dict(),
576             "scheduler_state_dict": scheduler.state_dict(),
577             "best_dev_dice": best_dev_dice,
578             "best_dev_miou": best_dev_miou,
579             "best_dev_loss": best_dev_loss,
580             # "history": agg,
581             "notes": f"Model saved based on best Micro Dice. Ignored index for metric: {ignore_index}"
582         }
583         torch.save(checkpoint, f"{model_save_dir}/{model_save_name}")
584         # Save model weights only
585         checkpoint = {
586             "epoch": t + 1,
587             "model_state_dict": model.state_dict(),
588         }
589         torch.save(checkpoint, f"{model_save_dir}/MO_{model_save_name}")
590     elif save:
591         print(f"Validation IoU score improved ({best_dev_miou:.6f}). Saving model...")
592
593         checkpoint = {
594             "epoch": t + 1,
595             "model_state_dict": model.state_dict(),
596             "optimizer_state_dict": optimizer.state_dict(),
597             "best_dev_dice": best_dev_dice,
598             "best_dev_miou": best_dev_miou,
599             "best_dev_loss": best_dev_loss,
600             # "history": agg,
601             "notes": f"Model saved based on best Micro Dice. Ignored index for metric: {ignore_index}"
602         }
603         torch.save(checkpoint, f"{model_save_dir}/{model_save_name}")
604         # Save model weights only
605         checkpoint = {
606             "epoch": t + 1,
607             "model_state_dict": model.state_dict(),
608         }
609         torch.save(checkpoint, f"{model_save_dir}/MO_{model_save_name}")
610     else:
611         print(f"Validation IoU score did not improve from {best_dev_miou:.6f}")
612
613
614     print("\n--- Training Finished! ---")
615     print(f"Best validation IoU score achieved: {best_dev_miou:.6f}")
616     print(f"Corresponding validation dice: {best_dev_dice:.6f}")
617     print(f"Corresponding validation loss: {best_dev_loss:.6f}")
618     print(f"Best model saved to: {os.path.join(model_save_dir, model_save_name)}")

```

Utils (utils/utils.py)

```

1 import torch
2 import matplotlib.pyplot as plt
3 import numpy as np
4 import torchvision.transforms.functional
5 import torch.nn.functional as F

```

```

6 from torchvision.transforms import InterpolationMode
7 from torchvision.transforms import functional as TF
8 from tqdm import tqdm
9 from PIL import Image
10 from typing import Optional, List
11
12
13 def resize_with_padding(image, target_size=512, interpolation=InterpolationMode.BILINEAR):
14     """
15     Resize a single image (Tensor of shape (C, H, W)) so that the longer side
16     equals target_size, preserving aspect ratio; add black padding as needed.
17     Args:
18         image (Tensor): Input image tensor of shape (C, H, W).
19         target_size (int): The target size for the longer side of the image.
20         interpolation (InterpolationMode): The interpolation method to use.
21     Returns:
22         tuple: A tuple containing the resized and padded image, plus a metadata dictionary.
23     """
24     _, orig_h, orig_w = image.shape
25     scale = min(target_size / orig_w, target_size / orig_h)
26     new_w = int(round(orig_w * scale))
27     new_h = int(round(orig_h * scale))
28
29     # Resize the image
30     image_resized = TF.resize(image, size=(new_h, new_w), interpolation=interpolation)
31
32     # Compute padding on each side
33     pad_w = target_size - new_w
34     pad_h = target_size - new_h
35     pad_left = pad_w // 2
36     pad_right = pad_w - pad_left
37     pad_top = pad_h // 2
38     pad_bottom = pad_h - pad_top
39
40     # Pad the image (padding order: left, top, right, bottom)
41     image_padded = TF.pad(image_resized, padding=(pad_left, pad_top, pad_right, pad_bottom), fill=0)
42
43     meta = {
44         "original_size": (orig_h, orig_w),
45         "new_size": (new_h, new_w),
46         "pad": (pad_left, pad_top, pad_right, pad_bottom),
47         "scale": scale
48     }
49     return image_padded, meta
50
51 def reverse_resize_and_padding(image, meta, interpolation="bilinear"):
52     """
53     Remove the padding from image (Tensor of shape (C, target_size, target_size))
54     using metadata and then resize the cropped image back to the original size.
55     Args:
56         image (Tensor): The input image tensor.
57         meta (dict): The metadata dictionary containing information about the original and new sizes and padding.
58         interpolation (str): The interpolation method to use ("bilinear" or "nearest").
59     Returns:
60         Tensor: The image with padding removed and resized to the original size.
61     """
62     pad_left, pad_top, pad_right, pad_bottom = meta["pad"]
63     new_h, new_w = meta["new_size"]
64
65     # Crop out the padding: from pad_top to pad_top+new_h and pad_left to pad_left+new_w.
66     image_cropped = image[..., pad_top: pad_top + new_h, pad_left: pad_left + new_w]
67
68     # Resize the cropped image back to the original size.
69     orig_h, orig_w = meta["original_size"]
70     # F.interpolate expects a 4D tensor.
71     image_original = F.interpolate(image_cropped.unsqueeze(0),
72                                     size=(orig_h, orig_w),

```

```

73         mode=interpolation,
74         align_corners=False if interpolation != "nearest" else None)
75     return image_original.squeeze(0)
76
77 def process_batch_forward(batch_images, target_size=512, interpolation=InterpolationMode.BILINEAR):
78     """
79     Process a batch (Tensor of shape (N, C, H, W)) by resizing each image to target_size
80     with aspect ratio preserved (adding black padding).
81     Args:
82         batch_images (Tensor): A batch of images (N, C, H, W).
83         target_size (int): The target size for the longer side of the image.
84         interpolation (InterpolationMode): The interpolation method to use.
85     Returns:
86         tuple: A tuple containing the processed batch and a list of meta dictionaries.
87     """
88     resized_batch = []
89     meta_list = []
90     for image in batch_images:
91         # If the image has 4 channels, slice to keep only the first 3 (RGB).
92         if image.ndim == 3 and image.shape[0] == 4:
93             image = image[:3, ...]
94         image_resized, meta = resize_with_padding(image, target_size, interpolation)
95         resized_batch.append(image_resized)
96         meta_list.append(meta)
97     return torch.stack(resized_batch), meta_list
98
99 def process_batch_reverse(batch_outputs, meta_list, interpolation="bilinear"):
100     """
101     Given a batch of network outputs of shape (N, C, target_size, target_size) and the
102     corresponding meta info, reverse the transform for each one to obtain predictions at their
103     original sizes.
104     Args:
105         batch_outputs (Tensor): A batch of network outputs (N, C, target_size, target_size).
106         meta_list (list): A list of meta dictionaries, one for each image in the batch.
107         interpolation (str): The interpolation method to use ("bilinear" or "nearest").
108     Returns:
109         list: A list of tensors, each with the original size.
110     """
111     original_outputs = []
112     for output, meta in zip(batch_outputs, meta_list):
113         restored = reverse_resize_and_padding(output, meta, interpolation=interpolation)
114         original_outputs.append(restored)
115     return original_outputs
116
117 def calculate_class_weights(
118     label_source,
119     num_classes: int,
120     ignore_index: Optional[int] = None,
121     source_type: str = 'files',
122     unimportant_class_indices: Optional[List[int]] = None, # Indices to down-weight
123     target_unimportant_weight: float = 1.0, # Target weight for unimportant classes
124     normalize_target_sum: float = -1.0 # Normalize weights sum (-1 means num_classes)
125 ) -> torch.Tensor:
126     """
127     Calculates class weights based on inverse frequency, then adjusts weights
128     for specified unimportant classes and re-normalizes.
129     Args:
130         label_source: Source of labels, can be a list of file paths or a dataset.
131         num_classes (int): The total number of classes.
132         ignore_index (Optional[int]): Index to ignore in the labels.
133         source_type (str): Type of source, 'files' or 'dataset'.
134         unimportant_class_indices (Optional[List[int]]): Indices of unimportant classes.
135         target_unimportant_weight (float): Target weight for unimportant classes.
136         normalize_target_sum (float): Normalize weights sum. If -1, normalize to num_classes.
137     Returns:
138         torch.Tensor: A tensor of class weights.
139     """

```

```

140
141 counts = torch.zeros(num_classes, dtype=torch.float64)
142 total_valid_pixels = 0
143
144 iterator = None
145 num_labels = 0
146 if source_type == 'files':
147     iterator = label_source
148     num_labels = len(label_source)
149 elif source_type == 'dataset':
150     iterator = range(len(label_source))
151     num_labels = len(label_source)
152 else:
153     raise ValueError("source_type must be either 'files' or 'dataset'")
154
155 print(f"Processing {num_labels} labels...")
156 pbar = tqdm(iterator, total=num_labels)
157 for idx_or_path in pbar:
158     label_data = None
159     if source_type == 'files':
160         path = idx_or_path; img = Image.open(path)
161         label_data = torch.from_numpy(np.array(img))
162     elif source_type == 'dataset':
163         _, label_data = label_source[idx_or_path];
164         label_data = torch.tensor(label_data) if not isinstance(label_data, torch.Tensor) else label_data
165
166     label_long = label_data.long().view(-1)
167
168     if ignore_index is not None:
169         valid_mask = (label_long != ignore_index)
170         label_valid = label_long[valid_mask]
171     else:
172         label_valid = label_long
173     label_valid = torch.clamp(label_valid, 0, num_classes - 1)
174
175     if label_valid.numel() > 0:
176         counts += torch.bincount(label_valid, minlength=num_classes).double()
177         total_valid_pixels += label_valid.numel()
178
179 print("\nFinished counting.")
180 print(f"Raw pixel counts per class: {counts.long().tolist()}")
181 print(f"Total valid pixels counted: {total_valid_pixels}")
182
183 frequencies = counts / total_valid_pixels
184 epsilon = 1e-6
185 inverse_frequencies = 1.0 / (frequencies + epsilon)
186
187 weights = inverse_frequencies
188
189 if unimportant_class_indices:
190     for idx in unimportant_class_indices:
191         weights[idx] = min(weights)
192
193 target_sum = normalize_target_sum if normalize_target_sum > 0 else float(num_classes)
194 final_weights = weights / weights.sum() * target_sum
195
196 print(f"Calculated Final Class Weights: {final_weights.tolist()}")
197
198 return final_weights.float()
199
200
201 def convert_rgb_label_to_classes(label_array_rgb):
202     """
203     Converts a 3-channel RGB label map to a 1-channel class map.
204
205     Mapping:
206     [0, 0, 0] (Black)      -> 0 (Background)

```

```

207     [128, 0, 0] (Red)          -> 1 (Cat)
208     [0, 128, 0] (Green)      -> 2 (Dog)
209     [255, 255, 255] (White) -> 0 (Background) - Assuming white is also background
210     Other                    -> 255 (Ignore)
211
212     Args:
213         label_array_rgb (np.ndarray): A HxWx3 NumPy array (uint8).
214
215     Returns:
216         np.ndarray: A HxW NumPy array (uint8) with class indices.
217     """
218     # Input validation
219     if label_array_rgb.ndim != 3 or label_array_rgb.shape[2] != 3:
220         raise ValueError(
221             "Input label must be 3-channel RGB (HxWx3), "
222             f"but got shape {label_array_rgb.shape}"
223         )
224
225     h, w, _ = label_array_rgb.shape
226     # Initialize with ignore value (255)
227     label_map_1channel = np.full((h, w), 255, dtype=np.uint8)
228
229     # Define colors as tuples for comparison
230     black = (0, 0, 0)
231     red = (128, 0, 0)
232     green = (0, 128, 0)
233     white = (255, 255, 255)
234
235     # Create boolean masks for each color
236     # Comparing tuples is faster for multi-channel exact matches typically
237     mask_black = np.all(label_array_rgb == black, axis=2)
238     mask_red = np.all(label_array_rgb == red, axis=2)
239     mask_green = np.all(label_array_rgb == green, axis=2)
240     mask_white = np.all(label_array_rgb == white, axis=2)
241
242     # Apply mapping (order can matter if masks could overlap, but shouldn't here)
243     # Map backgrounds first
244     label_map_1channel[mask_black] = 0
245     label_map_1channel[mask_white] = 0 # Map white to background class 0
246     # Map foreground classes
247     label_map_1channel[mask_red] = 1 # Cat
248     label_map_1channel[mask_green] = 2 # Dog
249     # Any remaining pixels stay 255
250
251     return label_map_1channel

```

Losses (utils/weighted_loss.py)

```

1  import torch
2  import torch.nn as nn
3  import torch.nn.functional as F
4  from typing import Optional, Callable
5
6  class WeightedMemoryEfficientDiceLoss(nn.Module):
7      """
8      Calculates a memory-efficient Soft Dice loss, optionally with class weights.
9      Args:
10         apply_softmax (bool): Whether to apply softmax to the input logits. Defaults to True.
11         ignore_index (int, optional): Index to ignore in the loss calculation. Defaults to None.
12         class_weights (torch.Tensor, optional): Weights for each class. Defaults to None.
13         smooth (float): Smoothing factor to prevent division by zero. Defaults to 1e-5.
14      """
15     def __init__(self,
16                 apply_softmax: bool = True,
17                 ignore_index: Optional[int] = None,
18                 class_weights: Optional[torch.Tensor] = None,

```



```

19         smooth: float = 1e-5):
20     super().__init__()
21     self.apply_softmax = apply_softmax
22     self.ignore_index = ignore_index
23     self.smooth = smooth
24
25     if class_weights is not None:
26         self.class_weights = class_weights
27     else:
28         self.class_weights = None
29
30
31     def forward(self, x, y):
32
33         num_classes = x.shape[1]
34         shp_y = y.shape
35         if self.apply_softmax:
36             probs = F.softmax(x, dim=1)
37         else:
38             probs = x
39
40         with torch.no_grad():
41             # Handle potential shape mismatches between predictions and targets
42             if len(shp_y) != len(probs.shape):
43                 if len(shp_y) == len(probs.shape) - 1 and len(shp_y) >= 2 and shp_y == probs.shape[2:]:
44                     y = y.unsqueeze(1)
45                 elif len(shp_y) == len(probs.shape) and shp_y[1] == 1: pass # ok
46                 else: raise ValueError(f"Shape mismatch: probs {probs.shape}, y {shp_y}")
47             y_long = y.long()
48
49             mask = None
50             if probs.shape == y.shape:
51                 y_onehot = y.float()
52                 if mask is not None:
53                     y_indices_for_mask = torch.argmax(y_onehot, dim=1, keepdim=True)
54                     mask = (y_indices_for_mask != self.ignore_index)
55                     y_onehot = y_onehot * mask
56             else:
57                 y_onehot = torch.zeros_like(probs, device=probs.device)
58                 y_onehot.scatter_(1, y_long, 1)
59                 if mask is not None: y_onehot = y_onehot * mask
60
61             sum_gt = y_onehot.sum(dim=(2, 3))
62
63             if mask is not None:
64                 probs = probs * mask
65
66             intersect_persample = (probs * y_onehot).sum(dim=(2, 3))
67             sum_pred_persample = probs.sum(dim=(2, 3))
68             sum_gt_persample = sum_gt
69
70             # Aggregate across batch
71             intersect = intersect_persample.sum(0)
72             sum_pred = sum_pred_persample.sum(0)
73             sum_gt = sum_gt_persample.sum(0)
74
75             # Dice
76             denominator = sum_pred + sum_gt
77             dc = (2. * intersect + self.smooth) / (torch.clip(denominator + self.smooth, 1e-8))
78
79             valid_classes_mask = torch.ones_like(dc, dtype=torch.bool)
80             if self.ignore_index is not None and 0 <= self.ignore_index < num_classes:
81                 valid_classes_mask[self.ignore_index] = False
82
83             dc_final = torch.tensor(0.0, device=dc.device)
84
85             dc_valid = dc[valid_classes_mask]

```

```

86
87     if self.class_weights is not None:
88         # Use weighted average for valid classes
89         weights = self.class_weights.to(dc_valid.device)
90         weights_valid = weights[valid_classes_mask]
91
92         weighted_sum = (dc_valid * weights_valid).sum()
93         weight_sum = weights_valid.sum()
94         dc_final = weighted_sum / weight_sum.clamp(min=1e-8)
95     else:
96         dc_final = dc_valid.mean()
97
98     return -dc_final # Return negative Dice score as loss
99
100
101 class WeightedDiceCELoss(nn.Module):
102     """
103     Combines WeightedMemoryEfficientDiceLoss and Cross Entropy Loss with class_weights support.
104     Args:
105         dice_weight (float): Weight for the Dice loss. Defaults to 1.0.
106         ce_weight (float): Weight for the Cross Entropy loss. Defaults to 1.0.
107         ignore_index (int, optional): Index to ignore in the loss calculation. Defaults to None.
108         class_weights (torch.Tensor, optional): Weights for each class. Defaults to None.
109         smooth_dice (float): Smoothing factor for Dice loss. Defaults to 1e-5.
110         ce_kwargs (dict): Keyword arguments for CrossEntropyLoss. Defaults to {}.
111     """
112     def __init__(self,
113                 dice_weight: float = 1.0,
114                 ce_weight: float = 1.0,
115                 ignore_index: Optional[int] = None,
116                 class_weights: Optional[torch.Tensor] = None,
117                 smooth_dice: float = 1e-5,
118                 ce_kwargs={}):
119         super().__init__()
120         self.dice_weight = dice_weight
121         self.ce_weight = ce_weight
122         self.ignore_index = ignore_index
123
124         self.dice = WeightedMemoryEfficientDiceLoss(
125             apply_softmax=True,
126             ignore_index=ignore_index,
127             class_weights=class_weights,
128             smooth=smooth_dice
129         )
130
131         ce_final_kwargs = ce_kwargs.copy()
132         if ignore_index is not None:
133             ce_final_kwargs['ignore_index'] = ignore_index
134         if class_weights is not None:
135             ce_final_kwargs['weight'] = class_weights
136
137         self.cross_entropy = nn.CrossEntropyLoss(**ce_final_kwargs)
138
139     def forward(self, outputs, targets):
140         if targets.ndim == 3:
141             targets_dice = targets.unsqueeze(1).long() # Assuming [N, H, W]
142         elif targets.ndim == 4 and targets.shape[1] == 1:
143             targets_dice = targets.long()
144         else:
145             if targets.ndim == outputs.ndim and targets.shape[1] != 1:
146                 raise ValueError(f"Target shape {targets.shape} has multiple channels but expected class indices [N, H, W]")
147             else:
148                 raise ValueError(f"Unsupported target shape {targets.shape} for CE. Expected [N, H, W] or [N, 1, H, W]")
149
150         dice_loss = self.dice(outputs, targets_dice)
151
152

```

```

153     if targets.ndim == 4 and targets.shape[1] == 1:
154         targets_ce = targets.squeeze(1).long()
155     elif targets.ndim == 3:
156         targets_ce = targets.long()
157     else:
158         if targets.ndim == outputs.ndim and targets.shape[1] != 1:
159             raise ValueError(f"Target shape {targets.shape} has multiple channels but expected class indices [N, H, W]")
160         else:
161             raise ValueError(f"Unsupported target shape {targets.shape} for CE. Expected [N, H, W] or [N, 1, H, W]")
162
163     ce_loss = self.cross_entropy(outputs, targets_ce)
164
165     combined_loss = (self.dice_weight * dice_loss) + (self.ce_weight * ce_loss)
166     return combined_loss
167
168
169
170 class WeightedMemoryEfficientDiceLossPrompt(nn.Module):
171     """
172     Calculates a memory-efficient Dice loss, optionally with class weights and a non-linearity.
173     Args:
174         dice_nonlin (Callable, optional): Non-linearity function to apply to the predictions. Defaults to None.
175         apply_softmax (bool): Whether to apply softmax to the input logits. Defaults to True.
176         ignore_index (int, optional): Index to ignore in the loss calculation. Defaults to None.
177         class_weights (torch.Tensor, optional): Weights for each class. Defaults to None.
178         smooth (float): Smoothing factor to prevent division by zero. Defaults to 1e-5.
179     """
180     def __init__(self,
181                 dice_nonlin: Callable = None,
182                 apply_softmax: bool = True,
183                 ignore_index: Optional[int] = None,
184                 class_weights: Optional[torch.Tensor] = None, # New parameter
185                 smooth: float = 1e-5):
186         super().__init__()
187         self.apply_softmax = apply_softmax
188         self.ignore_index = ignore_index
189         self.smooth = smooth
190         self.dice_nonlin = dice_nonlin
191
192         # Store class weights, ensuring they are a Tensor if provided
193         if class_weights is not None:
194             assert isinstance(class_weights, torch.Tensor), "class_weights must be a torch.Tensor"
195             self.class_weights = class_weights
196         else:
197             self.class_weights = None
198
199
200     def forward(self, x, y):
201         num_classes = x.shape[1]
202         shp_y = y.shape
203
204         if self.apply_softmax:
205             probs = F.softmax(x, dim=1)
206         else:
207             probs = x
208
209         if self.dice_nonlin is not None:
210             probs = self.dice_nonlin(probs)
211
212         with torch.no_grad():
213             if len(shp_y) != len(probs.shape):
214                 if len(shp_y) == len(probs.shape) - 1 and len(shp_y) >= 2 and shp_y == probs.shape[2:]:
215                     y = y.unsqueeze(1)
216                 elif len(shp_y) == len(probs.shape) and shp_y[1] == 1: pass # ok
217                 else: raise ValueError(f"Shape mismatch: probs {probs.shape}, y {shp_y}")
218             y_long = y.long()
219

```

```

220     mask = None
221     if probs.shape == y.shape:
222         y_onehot = y.float()
223         if mask is not None:
224             y_indices_for_mask = torch.argmax(y_onehot, dim=1, keepdim=True)
225             mask = (y_indices_for_mask != self.ignore_index)
226             y_onehot = y_onehot * mask
227     else:
228         y_onehot = torch.zeros_like(probs, device=probs.device)
229         y_onehot.scatter_(1, y_long, 1)
230         if mask is not None: y_onehot = y_onehot * mask
231
232     sum_gt = y_onehot.sum(dim=(2, 3))
233
234     if mask is not None:
235         probs = probs * mask
236
237     intersect_persample = (probs * y_onehot).sum(dim=(2, 3))
238     sum_pred_persample = probs.sum(dim=(2, 3))
239     sum_gt_persample = sum_gt
240
241     # Aggregate across batch
242     intersect = intersect_persample.sum(0)
243     sum_pred = sum_pred_persample.sum(0)
244     sum_gt = sum_gt_persample.sum(0)
245
246     # Dice
247     denominator = sum_pred + sum_gt
248     dc = (2. * intersect + self.smooth) / (torch.clip(denominator + self.smooth, 1e-8))
249
250     valid_classes_mask = torch.ones_like(dc, dtype=torch.bool)
251     if self.ignore_index is not None and 0 <= self.ignore_index < num_classes:
252         valid_classes_mask[self.ignore_index] = False
253
254     dc_final = torch.tensor(0.0, device=dc.device)
255     dc_valid = dc[valid_classes_mask]
256     if self.class_weights is not None:
257         weights = self.class_weights.to(dc_valid.device)
258         weights_valid = weights[valid_classes_mask]
259         weighted_sum = (dc_valid * weights_valid).sum()
260         weight_sum = weights_valid.sum()
261         dc_final = weighted_sum / weight_sum.clamp(min=1e-8)
262     else:
263         dc_final = dc_valid.mean()
264
265     return -dc_final # Return negative Dice score as loss
266
267
268 class WeightedDiceNLLLoss(nn.Module):
269     """
270     Combines WeightedMemoryEfficientDiceLossPrompt and Cross Entropy Loss with class_weights support.
271     Args:
272         dice_weight (float): Weight for the Dice loss. Defaults to 1.0.
273         nll_weight (float): Weight for the NLL loss. Defaults to 1.0.
274         ignore_index (int, optional): Index to ignore in the loss calculation. Defaults to None.
275         class_weights (torch.Tensor, optional): Weights for each class. Defaults to None.
276         smooth_dice (float): Smoothing factor for Dice loss. Defaults to 1e-5.
277         apply_softmax (bool): Whether to apply softmax to the input logits. Defaults to True.
278         dice_nonlin (Callable, optional): Non-linearity function to apply to the predictions for Dice loss. Defaults to None.
279         nll_nonlin (Callable, optional): Non-linearity function to apply to the predictions for NLL loss. Defaults to None.
280         nll_kwargs (dict): Keyword arguments for NLLLoss. Defaults to {}.
281     """
282     def __init__(self,
283                 dice_weight: float = 1.0,
284                 nll_weight: float = 1.0,
285                 ignore_index: Optional[int] = None,
286                 class_weights: Optional[torch.Tensor] = None,

```

```

287         smooth_dice: float = 1e-5,
288         apply_softmax: bool = True,
289         dice_nonlin: Callable = None,
290         nll_nonlin: Callable = None,
291         nll_kwargs={}):
292     super().__init__()
293     self.dice_weight = dice_weight
294     self.nll_weight = nll_weight
295     self.ignore_index = ignore_index
296     self.dice_nonlin = dice_nonlin
297     self.nll_nonlin = nll_nonlin
298
299     self.dice = WeightedMemoryEfficientDiceLossPrompt(
300         apply_softmax=apply_softmax,
301         ignore_index=ignore_index,
302         class_weights=class_weights,
303         smooth=smooth_dice
304     )
305
306     nll_final_kwargs = nll_kwargs.copy()
307     if ignore_index is not None:
308         nll_final_kwargs['ignore_index'] = ignore_index
309     if class_weights is not None:
310         nll_final_kwargs['weight'] = class_weights
311
312     self.nll = nn.NLLLoss(**nll_final_kwargs)
313
314     def forward(self, outputs, targets):
315         if targets.ndim == 3:
316             targets_dice = targets.unsqueeze(1).long()
317         elif targets.ndim == 4 and targets.shape[1] == 1:
318             targets_dice = targets.long()
319         else:
320             if targets.ndim == outputs.ndim and targets.shape[1] != 1:
321                 raise ValueError(f"Target shape {targets.shape} has multiple channels but expected class indices [N, H, W] or [N, 1, H, W]")
322             else:
323                 raise ValueError(f"Unsupported target shape {targets.shape} for CE. Expected [N, H, W] or [N, 1, H, W]")
324
325         dice_loss = self.dice(outputs, targets_dice)
326
327         if targets.ndim == 4 and targets.shape[1] == 1:
328             targets_nll = targets.squeeze(1).long()
329         elif targets.ndim == 3:
330             targets_nll = targets.long()
331         else:
332             if targets.ndim == outputs.ndim and targets.shape[1] != 1:
333                 raise ValueError(f"Target shape {targets.shape} has multiple channels but expected class indices [N, H, W] or [N, 1, H, W]")
334             else:
335                 raise ValueError(f"Unsupported target shape {targets.shape} for CE. Expected [N, H, W] or [N, 1, H, W]")
336
337         if self.nll_nonlin is not None:
338             outputs = self.nll_nonlin(outputs)
339         nll_loss = self.nll(outputs, targets_nll)
340
341         combined_loss = (self.dice_weight * dice_loss) + (self.nll_weight * nll_loss)
342         return combined_loss

```

Augmentations (utils/augmentation.ipynb)

Resize

```

1 import imgaug.augmenters as iaa
2 from imgaug.augmentables.segmaps import SegmentationMapsOnImage
3
4 # Padding augmenter

```

```

5 pad_aug = iaa.PadToAspectRatio(
6     1.0, # Keep aspect ratio
7     position="center",
8     pad_mode="constant",
9     pad_cval=0
10 )
11
12 # Resizers for Image and Labels
13 resize_img = iaa.Resize(256, interpolation="cubic")
14 resize_mask = iaa.Resize(256, interpolation="nearest")
15
16 # Function to resize images
17 def image_resizer(images, random_state, parents, hooks):
18     """
19     Resizes the images using the resize_img augmenter.
20     Args:
21         images (list): A list of images to resize.
22         random_state: Random state for augmentation.
23         parents: Parents for augmentation.
24         hooks: Hooks for augmentation.
25     Returns:
26         list: A list of resized images.
27     """
28     return [resize_img.augment_image(img) for img in images]
29
30 # Function to resize masks
31 def label_resizer(segmaps, random_state, parents, hooks):
32     """
33     Resizes the segmentation masks using the resize_mask augmenter.
34     Args:
35         segmaps (list): A list of SegmentationMapsOnImage objects to resize.
36         random_state: Random state for augmentation.
37         parents: Parents for augmentation.
38         hooks: Hooks for augmentation.
39     Returns:
40         list: A list of resized SegmentationMapsOnImage objects.
41     """
42     new_segmaps = []
43     for segmap in segmaps:
44         new_arr = resize_mask.augment_image(segmap.arr) # Resize the mask array.
45         new_segmaps.append(SegmentationMapsOnImage(new_arr, shape=new_arr.shape)) # Create a new SegmentationMapsOnImage
46     return new_segmaps
47
48 # Resizing augmenter
49 resize_aug = iaa.Sequential([
50     pad_aug, # Apply padding to maintain aspect ratio.
51     iaa.Lambda( # Apply lambda to resize images and masks.
52         func_images=image_resizer, # Function to resize images.
53         func_segmentation_maps=label_resizer, # Function to resize segmentation masks.
54     )
55 ])

```

Rotation

```

1 import imgaug.augmenters as iaa
2 from imgaug.augmentables.segmaps import SegmentationMapsOnImage
3
4 # Applies rotation augmentation to both images and segmentation maps
5 rotation_aug = iaa.Sequential([
6     iaa.Affine(
7         rotate=(45, 315), # Rotation angle
8         fit_output=True, # Maintain full rotated image
9         mode="constant",
10         cval=0,
11         backend="cv2"

```

```

12     ),
13     resize_aug
14 ])

```

Random Cropping

```

1  import imgaug.augmenters as iaa
2  from imgaug.imgaug import SegmentationMapsOnImage
3
4  class CenterSquareCropAugmenter(iaa.Augmenter):
5      """
6      CenterSquareCropAugmenter crops images to a square shape, centered.
7      Args:
8          name (str, optional): Name of the augmenter. Defaults to None.
9          deterministic (bool, optional): Whether the augmentation is deterministic. Defaults to False.
10         random_state (None, optional): Random state. Defaults to None.
11     """
12     def __init__(self, name=None, deterministic=False, random_state=None):
13         super(CenterSquareCropAugmenter, self).__init__(
14             name=name, deterministic=deterministic, random_state=random_state)
15         self.cropper = iaa.CropToAspectRatio(1.0, position="center") # Square cropper
16
17     def _augment_images(self, images, random_state, parents, hooks):
18         """
19         Applies the center square crop to a list of images.
20         Args:
21             images (list of numpy.ndarray): List of images.
22             random_state (numpy.random.RandomState): Random state.
23             parents (imgaug.parameters.StochasticParameter): Parents.
24             hooks (imgaug.hook.HooksImages): Hooks.
25         Returns:
26             list of numpy.ndarray: List of cropped images.
27         """
28         return [self.cropper.augment_image(img) for img in images]
29
30     def _augment_segmentation_maps(self, segmaps, random_state, parents, hooks):
31         """
32         Applies the center square crop to segmentation maps.
33         Args:
34             segmaps (list of imgaug.imgaug.SegmentationMapsOnImage): List of segmentation maps.
35             random_state (numpy.random.RandomState): Random state.
36             parents (imgaug.parameters.StochasticParameter): Parents.
37             hooks (imgaug.hook.HooksImages): Hooks.
38         Returns:
39             list of imgaug.imgaug.SegmentationMapsOnImage: List of cropped segmentation maps.
40         """
41         out = []
42         for segmap in segmaps:
43             cropped_arr = self.cropper.augment_image(segmap.arr)
44             out.append(SegmentationMapsOnImage(cropped_arr, shape=cropped_arr.shape))
45         return out
46
47     def get_parameters(self):
48         return []
49
50
51     class RandomSquareCropAugmenter(iaa.Augmenter):
52         """
53         RandomSquareCropAugmenter crops images to a square shape, randomly.
54         Args:
55             crop_factor (float, optional): Ratio of the smallest edge to use for the square. Defaults to 2/3.
56             name (str, optional): Name of the augmenter. Defaults to None.
57             deterministic (bool, optional): Whether the augmentation is deterministic. Defaults to False.
58             random_state (None, optional): Random state. Defaults to None.
59         """

```



```

60 def __init__(self, crop_factor=2/3, name=None, deterministic=False, random_state=None):
61     """
62     crop_factor: Ratio of the smallest edge to use for the square.
63     """
64     super(RandomSquareCropAugmenter, self).__init__(
65         name=name, deterministic=deterministic, random_state=random_state)
66     self.crop_factor = crop_factor
67
68 def _augment_images(self, images, random_state, parents, hooks):
69     """
70     Applies the random square crop to a list of images.
71     Args:
72         images (list of numpy.ndarray): List of images.
73         random_state (numpy.random.RandomState): Random state.
74         parents (imgaug.parameters.StochasticParameter): Parents.
75         hooks (imgaug.hook.HooksImages): Hooks.
76     Returns:
77         list of numpy.ndarray: List of cropped images.
78     """
79     out_images = []
80     for img in images:
81         H, W = img.shape[:2]
82         min_side = min(H, W)
83         crop_size = int(min_side * self.crop_factor)
84         max_x = W - crop_size
85         max_y = H - crop_size
86         x1 = random_state.randint(0, max_x + 1)
87         y1 = random_state.randint(0, max_y + 1)
88         cropped_img = img[y1: y1 + crop_size, x1: x1 + crop_size]
89         out_images.append(cropped_img)
90     return out_images
91
92 def _augment_segmentation_maps(self, segmaps, random_state, parents, hooks):
93     """
94     Applies the random square crop to segmentation maps.
95     Args:
96         segmaps (list of imgaug.imgaug.SegmentationMapsOnImage): List of segmentation maps.
97         random_state (numpy.random.RandomState): Random state.
98         parents (imgaug.parameters.StochasticParameter): Parents.
99         hooks (imgaug.hook.HooksImages): Hooks.
100     Returns:
101         list of imgaug.imgaug.SegmentationMapsOnImage: List of cropped segmentation maps.
102     """
103     out_segmaps = []
104     for segmap in segmaps:
105         arr = segmap.arr
106         H, W = arr.shape[:2]
107         min_side = min(H, W)
108         crop_size = int(min_side * self.crop_factor)
109         max_x = W - crop_size
110         max_y = H - crop_size
111         x1 = random_state.randint(0, max_x + 1)
112         y1 = random_state.randint(0, max_y + 1)
113         cropped_arr = arr[y1: y1 + crop_size, x1: x1 + crop_size]
114         out_segmaps.append(SegmentationMapsOnImage(cropped_arr, shape=cropped_arr.shape))
115     return out_segmaps
116
117 def get_parameters(self):
118     return [self.crop_factor]
119
120 # Center crop augmenter.
121 center_crop_aug = iaa.Sequential([CenterSquareCropAugmenter(), resize_aug])
122
123 # Random square crop augmenter.
124 random_crop_aug = iaa.Sequential([RandomSquareCropAugmenter(), resize_aug])

```

Random Masking

```
1 import numpy as np
2 import imgaug.augmenters as iaa
3 from imgaug.augmentables.segmaps import SegmentationMapsOnImage
4
5 # Define augmentation for masking images.
6 mask_im_aug = iaa.Sequential([
7     iaa.CoarseDropout(p=0.15, size_percent=(1/50), random_state=2)
8 ])
9
10 # Define augmentation for masking labels.
11 mask_label_aug = iaa.Sequential([
12     iaa.CoarseDropout(p=0.15, size_percent=(1/50), random_state=2)
13 ])
14
15 def random_masking_labels(segmaps, random_state, parents, hooks):
16     """
17     Applies random masking to segmentation maps.
18     Args:
19         segmaps: Input segmentation maps.
20         random_state: Random state for augmentation.
21         parents: Parent objects.
22         hooks: Hooks for augmentation.
23     Returns:
24         List of augmented segmentation maps.
25     """
26     new_segmaps = []
27     for segmap in segmaps:
28         # Convert segmentation map array to uint8 for augmentation.
29         segmap_arr_uint8 = segmap.arr.astype(np.uint8)
30         # Apply mask_label_aug to the segmentation map.
31         new_arr = mask_label_aug.augment_image(segmap_arr_uint8)
32         # Create a new SegmentationMapsOnImage object with the augmented array.
33         new_segmaps.append(SegmentationMapsOnImage(new_arr, shape=new_arr.shape))
34     return new_segmaps
35
36 # Random Masking augmentation
37 masking_aug = iaa.Sequential([
38     iaa.Lambda(
39         func_images=lambda images, rs, parents, hooks: [mask_im_aug.augment_image(img) for img in images],
40         func_segmentation_maps=random_masking_labels
41     ),
42     resize_aug # Apply resizing after masking
43 ])
```

Grayscale

```
1 # Grayscale augmentation
2 grayscale = iaa.Grayscale(alpha=1.0, from_colorspace="RGB")
3 grayscale_aug = iaa.Sequential([
4     grayscale,
5     resize_aug
6 ])
```

Laplace Noise

```
1 # Laplace Noise Augmentation
2 laplace = iaa.AdditiveLaplaceNoise(scale=(0.1*255, 0.3*255), per_channel=True)
3 laplace_aug = iaa.Sequential([
4     laplace,
```

```
5     resize_aug
6 ])
```

Blur

```
1 # Blur Augmentation
2 blur = iaa.AverageBlur(k=(12))
3 blur_aug = iaa.Sequential([
4     blur,
5     resize_aug
6 ])
```

Contrast

```
1 # Contrast Augmentation
2 contrast = iaa.LinearContrast((0.2, 0.6))
3 contrast_aug = iaa.Sequential([
4     contrast,
5     resize_aug
6 ])
```

Merge

```
1 from PIL import Image
2 import math
3 import os
4 import numpy as np
5 from utils import convert_rgb_label_to_classes
6
7
8 def combine_images_preserve_aspect_ratio(image1_path, image2_path, output_path=None, is_label=False):
9     """
10     Combines two images, preserving aspect ratio, centers on 256x256, then optionally converts to a 1-channel class map
11     If is_label is True, converts the final RGB image to a 1-channel class map
12     using convert_rgb_label_to_classes before saving/returning.
13
14     Args:
15         image1_path (str): Path to the first image.
16         image2_path (str): Path to the second image.
17         output_path (str, optional): Path to save the final image. Defaults to None.
18         is_label (bool, optional): If True, apply label conversion. Defaults to False.
19
20     Returns:
21         PIL.Image.Image: The final combined image (RGB or L mode).
22
23     Raises:
24         FileNotFoundError, ValueError, IOError, RuntimeError as before.
25     """
26     TARGET_DIMENSION = 256
27     RESAMPLE_METHOD = Image.Resampling.NEAREST
28
29     def load_image(path):
30         """
31         Loads an image from the given path and converts it to RGB mode.
32         Handles RGBA, LA, and P modes by converting them to RGB to avoid issues.
33
34         Args:
35             path (str): The path to the image file.
36
37         Returns:
```

```

38         PIL.Image.Image: The loaded image in RGB mode.
39         """
40         img = Image.open(path)
41         if img.mode == 'RGBA':
42             # Create a new RGB image with a black background and paste the image, masking the alpha channel
43             background = Image.new('RGB', img.size, (0, 0, 0))
44             background.paste(img, mask=img.split()[-1])
45             img = background
46         elif img.mode == 'LA':
47             # Convert LA to RGBA, create a new RGB image with a black background, and paste the image, masking the alpha
48             rgba_img = img.convert('RGBA')
49             background = Image.new('RGB', rgba_img.size, (0, 0, 0))
50             background.paste(rgba_img, mask=rgba_img.split()[-1])
51             img = background
52         elif img.mode == 'P':
53             # Convert P to RGBA, create a new RGB image with a black background, and paste the image, masking the alpha
54             rgba_img = img.convert('RGBA')
55             background = Image.new('RGB', rgba_img.size, (0, 0, 0))
56             background.paste(rgba_img, mask=rgba_img.split()[-1])
57             img = background
58         return img.convert('RGB')
59
60
61     img1 = load_image(image1_path)
62     img2 = load_image(image2_path)
63
64     # 2. Dimensions & Orientation
65     w1, h1 = img1.size
66     w2, h2 = img2.size
67
68     def get_orientation(w, h):
69         """
70         Determine the orientation of an image (portrait or landscape).
71
72         Args:
73             w (int): Width of the image.
74             h (int): Height of the image.
75
76         Returns:
77             str: 'portrait' if the image is portrait, 'landscape' otherwise.
78         """
79         return 'portrait' if h > w else 'landscape'
80
81     orientation1 = get_orientation(w1, h1)
82     orientation2 = get_orientation(w2, h2)
83
84     if orientation1 != orientation2:
85         print(f" Mismatched orientations ({orientation1} vs {orientation2}) for {os.path.basename(image1_path)}, {os.path.basename(image2_path)}")
86         return None
87     orientation = orientation1
88
89     # 3. Calculate Scale
90     if orientation == 'portrait':
91         total_original_major_dim = w1 + w2
92         if total_original_major_dim == 0:
93             return None
94         scale = TARGET_DIMENSION / total_original_major_dim
95     else: # landscape
96         total_original_major_dim = h1 + h2
97         if total_original_major_dim == 0:
98             return None
99         scale = TARGET_DIMENSION / total_original_major_dim
100
101     # 4. Calculate Scaled Dimensions
102     scaled_w1 = max(1, math.ceil(w1 * scale))
103     scaled_h1 = max(1, math.ceil(h1 * scale))
104     scaled_w2 = max(1, math.ceil(w2 * scale))

```

```

105     scaled_h2 = max(1, math.ceil(h2 * scale))
106
107     # 5. Adjust for Exact Fit
108     final_w1, final_h1 = scaled_w1, scaled_h1
109     final_w2, final_h2 = scaled_w2, scaled_h2
110     if orientation == 'portrait':
111         diff = (scaled_w1 + scaled_w2) - TARGET_DIMENSION
112         if diff > 0:
113             final_w1 -= diff if scaled_w1 >= scaled_w2 else 0
114             final_w2 -= diff if scaled_w2 > scaled_w1 else 0
115     else:
116         diff = (scaled_h1 + scaled_h2) - TARGET_DIMENSION
117         if diff > 0:
118             final_h1 -= diff if scaled_h1 >= scaled_h2 else 0
119             final_h2 -= diff if scaled_h2 > scaled_h1 else 0
120
121     final_w1, final_h1, final_w2, final_h2 = max(1, final_w1), max(1, final_h1), max(1, final_w2), max(1, final_h2)
122
123     # 6. Resize Images
124
125     img1_resized = img1.resize((final_w1, final_h1), RESAMPLE_METHOD)
126     img2_resized = img2.resize((final_w2, final_h2), RESAMPLE_METHOD)
127
128
129     # 7. Create Combined Strip
130     if orientation == 'portrait':
131         combined_w, combined_h = TARGET_DIMENSION, max(final_h1, final_h2)
132         combined = Image.new('RGB', (combined_w, combined_h), (0, 0, 0))
133         combined.paste(img1_resized, (0, 0))
134         combined.paste(img2_resized, (final_w1, 0))
135     else: # landscape
136         combined_w, combined_h = max(final_w1, final_w2), TARGET_DIMENSION
137         combined = Image.new('RGB', (combined_w, combined_h), (0, 0, 0))
138         combined.paste(img1_resized, (0, 0))
139         combined.paste(img2_resized, (0, final_h1))
140
141     # 8. Create Final Canvas & Center
142     final_img = Image.new('RGB', (TARGET_DIMENSION, TARGET_DIMENSION), (0, 0, 0))
143     paste_x = (TARGET_DIMENSION - combined.width) // 2
144     paste_y = (TARGET_DIMENSION - combined.height) // 2
145     final_img.paste(combined, (paste_x, paste_y))
146
147     # 9. Label Conversion
148     if is_label:
149         # Convert final PIL Image to NumPy array
150         final_img_np = np.array(final_img)
151         # Apply the RGB -> Class ID conversion
152         label_map_lchannel = convert_rgb_label_to_classes(final_img_np)
153         # Convert the 1-channel NumPy array back to a PIL Image (mode 'L')
154         final_img = Image.fromarray(label_map_lchannel, mode='L')
155
156
157     # 10. Save if output path is provided
158     if output_path:
159         final_img.save(output_path)
160
161     return final_img

```

Augmentation without merge

```

1  import os
2  import random
3  import imageio
4  import numpy as np
5  from imgaug.augmentables.segmaps import SegmentationMapsOnImage

```

```

6 import math
7
8 # Define a dictionary mapping augementer names to their corresponding functions
9 augementer_dict = {
10     "rotation": rotation_aug,
11     "center_crop": center_crop_aug,
12     "random_crop": random_crop_aug,
13     "masking": masking_aug,
14     "grayscale": grayscale_aug,
15     "laplace": laplace_aug,
16     "blur": blur_aug,
17     "contrast": contrast_aug
18 }
19
20 # Calculate the number of augmenters
21 num_augmenters = len(augementer_dict) # Count based on the dictionary
22
23 # Configuration
24 folder_path = "Train/color" # Path to the folder containing color images
25 label_folder_path = "Train/label" # Path to the folder containing label images
26 save_color_dir = "astrain/color" # Directory to save augmented color images
27 save_label_dir = "astrain/label" # Directory to save augmented label images
28 majority_aug_factor = 1.5 # Augmentation factor to balance the dataset
29
30 # Create the output directories if they don't exist
31 os.makedirs(save_color_dir, exist_ok=True)
32 os.makedirs(save_label_dir, exist_ok=True)
33
34 # File Discovery and Classification
35 print("Scanning for image files...")
36 filenames = [
37     f for f in os.listdir(folder_path)
38     if f.lower().endswith(('.jpg', '.png'))
39 ]
40 # This section assumes that the file names can be used to identify the species
41
42 def get_species(filename):
43     """
44     Extracts the species name from a filename.
45     Args:
46         filename (str): The name of the file.
47     Returns:
48         str: The species name.
49     """
50     base = os.path.splitext(filename)[0]
51     parts = base.rsplit('_', 1)
52     return parts[0] if len(parts) > 1 else base
53
54 # Define a set of cat species for classification
55 cat_species = {
56     "Russian_Blue", "Siamese", "Sphynx", "Maine_Coon", "Abyssinian",
57     "Bombay", "British_Shorthair", "Bengal", "Egyptian_Mau", "Persian",
58     "Ragdoll", "Birman"
59 }
60
61 # Initialize lists to store cat and dog filenames
62 cat_files = []
63 dog_files = []
64 for fname in filenames:
65     species = get_species(fname)
66     name_no_ext = os.path.splitext(fname)[0]
67     label_path_check = os.path.join(label_folder_path, name_no_ext + ".png")
68     if species in cat_species:
69         cat_files.append(name_no_ext)
70     else:
71         dog_files.append(name_no_ext)
72

```

```

73 # Get the number of cat and dog files
74 N_cat = len(cat_files)
75 N_dog = len(dog_files)
76
77 print(f"Initial counts: Cats = {N_cat}, Dogs = {N_dog}")
78
79 # Copy Original Files
80 print("Processing originals with resize augmentation...")
81 all_original_files = cat_files + dog_files
82 processed_count = 0
83
84 # Ensure destination directories exist
85 os.makedirs(save_color_dir, exist_ok=True)
86 os.makedirs(save_label_dir, exist_ok=True)
87
88 for fname in all_original_files:
89     orig_color_path = os.path.join(folder_path, fname + ".jpg")
90     orig_label_path = os.path.join(label_folder_path, fname + ".png")
91
92     # Define destination paths (using original base name)
93     dest_color_path = os.path.join(save_color_dir, fname + ".jpg")
94     dest_label_path = os.path.join(save_label_dir, fname + ".png")
95
96     # Read the input image and its label
97     img = imageio.v2.imread(orig_color_path)
98     label = imageio.v2.imread(orig_label_path)
99
100     # Create a segmentation map object
101     segmap = SegmentationMapsOnImage(label, shape=img.shape)
102
103     # Apply the resize augmentation to both image and label map
104     resized_img, resized_segmap = resize_aug(image=img, segmentation_maps=segmap)
105     resized_label = resized_segmap.get_arr()
106
107     if resized_img.ndim == 3 and resized_img.shape[2] == 4:
108         resized_img = resized_img[..., :3] # RGBA to RGB
109
110     resized_label = convert_rgb_label_to_classes(resized_label)
111
112     # Ensure correct data types before saving
113     resized_img = resized_img.astype(np.uint8)
114     resized_label = resized_label.astype(np.uint8)
115
116     # Save the processed (resized) images
117     imageio.imwrite(dest_color_path, resized_img)
118     imageio.imwrite(dest_label_path, resized_label)
119
120     processed_count += 1
121
122 print(f"Processed and saved {processed_count} original image/label pairs using resize_aug.")
123
124
125 # --- Calculate Augmentation Needs ---
126 if N_cat == N_dog:
127     target_final_count = round(N_dog * majority_aug_factor)
128 elif N_cat < N_dog:
129     target_final_count = round(N_dog * majority_aug_factor)
130 else: # N_dog < N_cat
131     target_final_count = round(N_cat * majority_aug_factor)
132
133 total_aug_cat_needed = max(0, target_final_count - N_cat)
134 total_aug_dog_needed = max(0, target_final_count - N_dog)
135
136 print(f"Target final count per class: {target_final_count}")
137 print(f"Total augmentations needed: Cats = {total_aug_cat_needed}, Dogs = {total_aug_dog_needed}")
138
139 num_cats_per_aug = math.ceil(total_aug_cat_needed / num_augmenters)

```



```

140 num_dogs_per_aug = math.ceil(total_aug_dog_needed / num_augmenters)
141
142 print(f"Will select approximately {num_cats_per_aug} cats and {num_dogs_per_aug} dogs per augmentor.")
143
144 num_selected_cats = 0
145 num_selected_dogs = 0
146
147 # Augmentation Loop
148 generated_aug_count = 0
149 for i, (aug_name, aug_object) in enumerate(augmenter_dict.items()):
150     # Randomly select files for augmentation
151     selected_cats = random.choices(cat_files, k=num_cats_per_aug) if N_cat > 0 and num_cats_per_aug > 0 else []
152     selected_dogs = random.choices(dog_files, k=num_dogs_per_aug) if N_dog > 0 and num_dogs_per_aug > 0 else []
153     selected_files = selected_cats + selected_dogs
154     num_selected_cats += len(selected_cats)
155     num_selected_dogs += len(selected_dogs)
156
157     print(f"\nUsing augmentor '{aug_name}' ({i+1}/{num_augmenters}): processing {len(selected_files)} images ({len(selected_files)} processed)")
158
159     processed_in_batch = 0
160     for fname in selected_files:
161         color_path = os.path.join(folder_path, fname + ".jpg")
162         label_path = os.path.join(label_folder_path, fname + ".png")
163
164         # Read images
165         img = imageio.v2.imread(color_path)
166         label = imageio.v2.imread(label_path)
167         segmap = SegmentationMapsOnImage(label, shape=img.shape)
168
169         # Apply the augmentation
170         augmented_img, augmented_segmap = aug_object(image=img, segmentation_maps=segmap)
171         augmented_label = augmented_segmap.get_arr()
172
173         if augmented_img.ndim == 3 and augmented_img.shape[2] == 4:
174             augmented_img = augmented_img[..., :3] # RGBA to RGB
175
176         augmented_label = convert_rgb_label_to_classes(augmented_label)
177
178         # Type casting
179         augmented_label = augmented_label.astype(np.uint8)
180         augmented_img = augmented_img.astype(np.uint8)
181
182         out_color_path = os.path.join(save_color_dir, f"{fname}_{aug_name}_{processed_in_batch}.jpg")
183         out_label_path = os.path.join(save_label_dir, f"{fname}_{aug_name}_{processed_in_batch}.png")
184
185         # Save the augmented images
186         imageio.imwrite(out_color_path, augmented_img)
187         imageio.imwrite(out_label_path, augmented_label)
188
189         generated_aug_count += 1
190         processed_in_batch += 1
191
192     print(f"Augmenter '{aug_name}' finished. Processed {processed_in_batch} images.")

```

Augmentation with merge

```

1 import os
2 import random
3
4 SOURCE_COLOR_DIR = "Train/color"
5 SOURCE_LABEL_DIR = "Train/label"
6 DEST_COLOR_DIR = "astrain/color"
7 DEST_LABEL_DIR = "astrain/label"
8 NUM_COMBINATIONS_PER_TYPE = 126
9

```

```

10 cat_species = {
11     "Russian_Blue", "Siamese", "Sphynx", "Maine_Coon", "Abyssinian",
12     "Bombay", "British_Shorthair", "Bengal", "Egyptian_Mau", "Persian",
13     "Ragdoll", "Birman"
14 }
15
16 def get_species(filename):
17     """
18     Extracts the species name from a filename.
19     Args:
20         filename (str): The name of the file.
21     Returns:
22         str: The species name or the base filename if no species is found.
23     """
24     base = os.path.splitext(filename)[0]
25     parts = base.rsplit('_', 1)
26     return parts[0] if len(parts) > 1 else base
27
28
29 # 1. Create destination directories
30 os.makedirs(DEST_COLOR_DIR, exist_ok=True)
31 os.makedirs(DEST_LABEL_DIR, exist_ok=True)
32
33 # 2. Scan source directory and classify files
34 all_files_in_color = [
35     f for f in os.listdir(SOURCE_COLOR_DIR)
36     if f.lower().endswith(('.jpg', '.png')) # Assuming color can be jpg or png
37 ]
38
39
40 cat_files = []
41 dog_files = []
42
43 for fname_ext in all_files_in_color:
44     fname_no_ext = os.path.splitext(fname_ext)[0]
45     label_path_check = os.path.join(SOURCE_LABEL_DIR, fname_no_ext + ".png")
46
47
48     species = get_species(fname_ext)
49     if species in cat_species:
50         cat_files.append(fname_no_ext)
51     else:
52         dog_files.append(fname_no_ext)
53
54 N_cat = len(cat_files)
55 N_dog = len(dog_files)
56
57 print(f"Found {N_cat} cat images with labels.")
58 print(f"Found {N_dog} dog images with labels.")
59
60 # Function to generate combinations for a specific type
61 def generate_combinations(combo_type, files1_list, files2_list, num_required, output_prefix):
62     """
63     Generates N combinations by selecting files from lists and calling combine_images.
64     Prints the source files used for each successful combination.
65
66     Args:
67         combo_type (str): Description (e.g., "1 Cat + 1 Dog")
68         files1_list (list): List of base filenames for the first image.
69         files2_list (list): List of base filenames for the second image.
70         num_required (int): Target number of successful combinations.
71         output_prefix (str): Prefix for output filenames (e.g., "cat_dog").
72     """
73
74     combinations_done = 0
75     attempts = 0
76     max_attempts = num_required * 10

```

```

77
78 generated_pairs = set()
79 file1_base, file2_base = None, None
80
81 while combinations_done < num_required and attempts < max_attempts:
82     attempts += 1
83
84     if files1_list is files2_list:
85         # If combining from the same list, sample 2 unique files.
86         if len(files1_list) < 2:
87             break
88         file1_base, file2_base = random.sample(files1_list, 2)
89     else:
90         # If combining from different lists, sample one from each.
91         if not files1_list or not files2_list:
92             break
93         file1_base = random.choice(files1_list)
94         file2_base = random.choice(files2_list)
95
96     pair_key = tuple(sorted((file1_base, file2_base)))
97     if pair_key in generated_pairs:
98         # Skip if this pair has already been generated.
99         continue
100
101     # Construct paths
102     img1_color_ext = ".jpg"
103     img2_color_ext = ".jpg"
104     img1_label_ext = ".png"
105     img2_label_ext = ".png"
106
107     img1_color_path = os.path.join(SOURCE_COLOR_DIR, file1_base + img1_color_ext)
108     img1_label_path = os.path.join(SOURCE_LABEL_DIR, file1_base + img1_label_ext)
109     img2_color_path = os.path.join(SOURCE_COLOR_DIR, file2_base + img2_color_ext)
110     img2_label_path = os.path.join(SOURCE_LABEL_DIR, file2_base + img2_label_ext)
111
112     # Define output paths
113     output_base_name = f"{output_prefix}_{combinations_done}"
114     output_color_path = os.path.join(DEST_COLOR_DIR, output_base_name + ".jpg")
115     output_label_path = os.path.join(DEST_LABEL_DIR, output_base_name + ".png")
116
117     # Combine color images
118     combined_color = combine_images_preserve_aspect_ratio(img1_color_path, img2_color_path, output_color_path)
119
120     # Combine label images
121     combined_label = combine_images_preserve_aspect_ratio(img1_label_path, img2_label_path, output_label_path, True)
122
123     print(f"\n Generated: {output_base_name}.jpg/.png using [{file1_base}{img1_color_ext}], {file2_base}{img2_color_ext}")
124
125     combinations_done += 1
126     generated_pairs.add(pair_key)
127
128
129 # 1. 1 Cat + 1 Dog
130 generate_combinations(
131     combo_type="1 Cat + 1 Dog",
132     files1_list=cat_files,
133     files2_list=dog_files,
134     num_required=NUM_COMBINATIONS_PER_TYPE,
135     output_prefix="cat_dog"
136 )
137
138 # 2. 2 Cats
139 generate_combinations(
140     combo_type="2 Cats",
141     files1_list=cat_files,
142     files2_list=cat_files,
143     num_required=NUM_COMBINATIONS_PER_TYPE,

```

```

144     output_prefix="cat_cat"
145 )
146
147 # 3. 2 Dogs
148 generate_combinations(
149     combo_type="2 Dogs",
150     files1_list=dog_files,
151     files2_list=dog_files,
152     num_required=NUM_COMBINATIONS_PER_TYPE,
153     output_prefix="dog_dog"
154 )
155
156 print("\n--- Combination process finished. ---")

```

Augmentations for prompt based segmentation

```

1  import numpy as np
2  import random
3  import torch
4  import numpy as np
5  import random
6  import os
7  import time
8  import shutil
9  from torchvision.io import read_image
10 from PIL import Image
11 from utils.dataset import target_remap
12
13 def create_gaussian_heatmap(size=(256, 256), sigma=3.0):
14     """
15     Creates a 2D heatmap array with a Gaussian spot centered at a random pixel.
16
17     Args:
18         size (tuple): The (height, width) dimensions of the heatmap array.
19         sigma (float): The standard deviation (spread) of the Gaussian function.
20                       Larger sigma means a wider, smoother spot.
21
22     Returns:
23         numpy.ndarray: A 2D numpy array representing the heatmap (values typically 0-1).
24         tuple: The (y, x) coordinates of the chosen center pixel.
25     """
26     height, width = size
27     if height <= 0 or width <= 0:
28         raise ValueError("Size dimensions must be positive integers.")
29     if sigma <= 0:
30         raise ValueError("Sigma must be positive.")
31
32     # 1. Create a black canvas (array of zeros)
33     heatmap = np.zeros((height, width), dtype=np.float32) # Use float for calculations
34
35     # 2. Pick a random center pixel
36     center_y = random.randint(0, height - 1)
37     center_x = random.randint(0, width - 1)
38     print(f"Selected center pixel (y, x): ({center_y}, {center_x})")
39
40     # 3. Create coordinate grids
41     y_coords, x_coords = np.indices((height, width))
42
43     # 4. Calculate the squared Euclidean distance from the center for each pixel
44     dist_sq = (x_coords - center_x)**2 + (y_coords - center_y)**2
45
46     # 5. Calculate the Gaussian function
47     heatmap = np.exp(-dist_sq / (2 * sigma**2))
48
49     return heatmap, (center_y, center_x)

```

```

50
51
52 # --- Helper function for the selection process ---
53 def select_dominant_class(heatmap, remapped_mask):
54     """
55     Selects the dominant class in a mask based on heatmap scores.
56
57     Args:
58         heatmap (numpy.ndarray): The heatmap array.
59         remapped_mask (numpy.ndarray): The remapped mask array.
60
61     Returns:
62         int: The selected class (0 if no class is dominant).
63         dict: A dictionary of class scores.
64     """
65     class_scores = {}
66     present_classes = np.unique(remapped_mask)
67     target_classes = present_classes[present_classes > 0] # Classes 1, 2, 3
68
69     if target_classes.size == 0: return 0, {}
70
71     for class_val in target_classes:
72         mask_pixels = (remapped_mask == class_val)
73         if np.any(mask_pixels):
74             score = np.sum(heatmap[mask_pixels])
75             class_scores[class_val] = score
76         else:
77             class_scores[class_val] = 0
78
79     if not class_scores or all(s < 1e-9 for s in class_scores.values()):
80         selected_class = 0
81     else:
82         selected_class = max(class_scores, key=class_scores.get)
83
84     return selected_class, class_scores
85
86
87 TRAIN_IMG_DIR = "astrain/color"
88 TRAIN_LBL_DIR = "astrain/label"
89
90 HEATMAP_SIGMA = 3.0
91 MAX_ATTEMPTS = 1000
92
93 PSTRAIN_BASE_DIR = "pstrain" # New base output directory
94 PSTRAIN_IMG_DIR = os.path.join(PSTRAIN_BASE_DIR, "color") # For COPIED original images
95 PSTRAIN_HEATMAP_DIR = os.path.join(PSTRAIN_BASE_DIR, "point_prompt") # For heatmap IMAGES
96 PSTRAIN_LABEL_DIR = os.path.join(PSTRAIN_BASE_DIR, "label") # For final label masks
97
98
99 start_time = time.time()
100
101 os.makedirs(PSTRAIN_IMG_DIR, exist_ok=True)
102 os.makedirs(PSTRAIN_HEATMAP_DIR, exist_ok=True)
103 os.makedirs(PSTRAIN_LABEL_DIR, exist_ok=True)
104 print(f"Reading original images from: {os.path.abspath(TRAIN_IMG_DIR)}")
105 print(f"Reading labels from: {os.path.abspath(TRAIN_LBL_DIR)}")
106 print("-" * 30)
107 print(f"Saving copied images to: {os.path.abspath(PSTRAIN_IMG_DIR)}")
108 print(f"Saving heatmap images to: {os.path.abspath(PSTRAIN_HEATMAP_DIR)}")
109 print(f"Saving final label masks to: {os.path.abspath(PSTRAIN_LABEL_DIR)}")
110
111
112 all_label_files = os.listdir(TRAIN_LBL_DIR)
113 label_files = sorted([f for f in all_label_files if f.lower().endswith('.png') and not f.startswith('.')])
114
115
116 # --- Loop through all found label files ---

```

```

117 processed_count = 0
118 skipped_count = 0
119 error_count = 0
120 img_not_found_count = 0
121
122 total_files = len(label_files)
123 print(f"\nStarting processing for {total_files} label files...")
124
125 for i, label_filename in enumerate(label_files):
126     img_name_base = os.path.splitext(label_filename)[0] # Get base name without extension
127     label_filepath = os.path.join(TRAIN_LBL_DIR, label_filename)
128
129     print(f"\nProcessing label {i+1}/{total_files}: {label_filename} (Base: {img_name_base})")
130
131     # Find the corresponding original image file
132     original_img_path = None
133     original_img_ext = None
134     try:
135         found = False
136         for img_file in os.listdir(TRAIN_IMG_DIR):
137             if os.path.splitext(img_file)[0] == img_name_base:
138                 original_img_path = os.path.join(TRAIN_IMG_DIR, img_file)
139                 original_img_ext = os.path.splitext(img_file)[1] # Get extension (e.g., '.jpg')
140                 print(f" Found corresponding image: {img_file}")
141                 found = True
142                 break
143         if not found:
144             print(f" Skipping: Could not find corresponding image file for base name '{img_name_base}' in {TRAIN_IMG_DIR}")
145             img_not_found_count += 1
146             skipped_count += 1
147             continue
148     except Exception as e:
149         print(f"!!! ERROR searching for image file for {label_filename}: {e}")
150         error_count += 1
151         continue
152
153     try:
154         # Load the label mask file
155         label_tensor_loaded = read_image(label_filepath)
156
157         # Handle channel issues (ensure single channel)
158         if label_tensor_loaded.shape[0] != 1:
159             if label_tensor_loaded.shape[0] == 3:
160                 label_tensor_loaded = label_tensor_loaded[0:1, :, :]
161                 print(f" Info: Label had 3 channels, took the first.")
162             else:
163                 print(f" Skipping: Label has unexpected shape {label_tensor_loaded.shape}, expected (1, H, W).")
164                 skipped_count += 1
165                 continue
166
167         # Apply the first remap (255 -> 3)
168         label_tensor_original = target_remap(label_tensor_loaded)
169
170         # Process the Loaded Mask ONCE per sample
171         label_squeezed = label_tensor_original.squeeze(0)
172         mask_post_remap1 = label_squeezed.numpy().astype(np.uint8)
173         mask_size = mask_post_remap1.shape
174
175         # Apply the SECOND remapping (Swap 3->0, Add 1) -> Final classes 1, 2, 3
176         mask_swapped = mask_post_remap1.copy()
177         mask_swapped[mask_post_remap1 == 3] = 0
178         remapped_mask_final = mask_swapped + 1
179         final_present_classes = np.unique(remapped_mask_final)
180         final_target_classes = final_present_classes[final_present_classes > 0]
181
182         # Check if finding two distinct classes is possible

```

```

184     if len(final_target_classes) < 2:
185         print(f" Skipping: Mask only contains {len(final_target_classes)} target class(es) {final_target_classes}")
186         skipped_count += 1
187         continue
188
189     # Loop to find TWO distinct class selections for this sample
190     selected_results = [] # List to store (selected_class, final_mask_array, heatmap_array)
191     attempts = 0
192     found_classes = set()
193
194     while len(selected_results) < 2 and attempts < MAX_ATTEMPTS:
195         attempts += 1
196         heatmap, center_coords = create_gaussian_heatmap(size=mask_size, sigma=HEATMAP_SIGMA)
197         current_selected_class, _ = select_dominant_class(heatmap, remapped_mask_final)
198
199         if current_selected_class > 0 and current_selected_class not in found_classes:
200             final_mask = np.zeros_like(remapped_mask_final, dtype=np.uint8)
201             final_mask[remapped_mask_final == current_selected_class] = current_selected_class
202             selected_results.append((current_selected_class, final_mask, heatmap))
203             found_classes.add(current_selected_class)
204             print(f" Attempt {attempts}: Found distinct class {current_selected_class} at {center_coords}")
205
206     if len(selected_results) == 2:
207         print(f" Successfully found two distinct classes.")
208         sel_cls_1, fin_msk_1, heatmap_1 = selected_results[0]
209         sel_cls_2, fin_msk_2, heatmap_2 = selected_results[1]
210
211         # --- Define final output filenames (consistent naming) ---
212         output_base_name_1 = f"{img_name_base}_1"
213         output_base_name_2 = f"{img_name_base}_2"
214
215         # Paths for triplet 1
216         output_img1_path = os.path.join(PSTRAIN_IMG_DIR, f"{output_base_name_1}{original_img_ext}")
217         output_heatmap1_path = os.path.join(PSTRAIN_HEATMAP_DIR, f"{output_base_name_1}.png")
218         output_label1_path = os.path.join(PSTRAIN_LABEL_DIR, f"{output_base_name_1}.png")
219
220         # Paths for triplet 2
221         output_img2_path = os.path.join(PSTRAIN_IMG_DIR, f"{output_base_name_2}{original_img_ext}")
222         output_heatmap2_path = os.path.join(PSTRAIN_HEATMAP_DIR, f"{output_base_name_2}.png")
223         output_label2_path = os.path.join(PSTRAIN_LABEL_DIR, f"{output_base_name_2}.png")
224
225         # Copy the original image twice
226         shutil.copy2(original_img_path, output_img1_path)
227         shutil.copy2(original_img_path, output_img2_path)
228         print(f" Copied original image to: {os.path.basename(output_img1_path)}")
229         print(f" Copied original image to: {os.path.basename(output_img2_path)}")
230
231         # Save Heatmaps as PNG Images (Scaled 0-255)
232         heatmap1_scaled = (heatmap_1 * 255).astype(np.uint8)
233         heatmap2_scaled = (heatmap_2 * 255).astype(np.uint8)
234         Image.fromarray(heatmap1_scaled, mode='L').save(output_heatmap1_path) # 'L' mode for grayscale
235         Image.fromarray(heatmap2_scaled, mode='L').save(output_heatmap2_path)
236         print(f" Saved Heatmap 1: {os.path.basename(output_heatmap1_path)}")
237         print(f" Saved Heatmap 2: {os.path.basename(output_heatmap2_path)}")
238
239         # Save Final Masks as PNG Images (already uint8)
240         Image.fromarray(fin_msk_1).save(output_label1_path)
241         Image.fromarray(fin_msk_2).save(output_label2_path)
242         print(f" Saved Label 1: {os.path.basename(output_label1_path)}")
243         print(f" Saved Label 2: {os.path.basename(output_label2_path)}")
244
245         processed_count += 1 # Count original files that yielded 2 outputs
246     else:

```

```

251         print(f" Skipping: Failed to find two distinct classes within {MAX_ATTEMPTS} attempts.")
252         skipped_count += 1
253
254     except FileNotFoundError:
255         print(f"!!! ERROR processing label file {label_filename}: File not found (unexpected).")
256         error_count += 1
257     except Exception as e:
258         print(f"!!! ERROR processing label file {label_filename}: {e}")
259         import traceback
260         traceback.print_exc()
261         error_count += 1
262         # Continue to the next file
263
264
265 # --- Final Summary ---
266 end_time = time.time()
267 total_time = end_time - start_time
268 print("\n" + "="*40)
269 print("Processing Complete.")
270 print(f"Output base directory: {os.path.abspath(PSTRAIN_BASE_DIR)}")
271 print("-" * 40)
272 print(f"Total label files found: {total_files}")
273 print(f"Successfully processed (2 triplets generated): {processed_count}")
274 print(f"Skipped (due to various reasons): {skipped_count}")
275 print(f" - Skipped because original image not found: {img_not_found_count}")
276 print(f" - Skipped (other reasons, e.g., too few classes): {skipped_count - img_not_found_count}")
277 print(f"Errors during processing: {error_count}")
278 print("-" * 40)
279 print(f"Total files created in '{os.path.basename(PSTRAIN_IMG_DIR)}': {processed_count * 2}")
280 print(f"Total files created in '{os.path.basename(PSTRAIN_HEATMAP_DIR)}': {processed_count * 2}")
281 print(f"Total files created in '{os.path.basename(PSTRAIN_LABEL_DIR)}': {processed_count * 2}")
282 print("-" * 40)
283 print(f"Total time: {total_time:.2f} seconds")
284 print("="*40)

```

UNet training (unet.ipynb)

```

1 import torch
2 from torch import nn
3 from torch import Tensor
4 from torch import optim
5 from torch.utils.data import DataLoader
6 from utils.training import start
7 from utils.MetricsHistory import MetricsHistory
8 from unet.unet import unet
9 from utils.weighted_loss import WeightedDiceCELoss
10 from utils.utils import calculate_class_weights
11 from utils.dataset import dataset, target_remap, diff_size_collate
12
13 EVAL_IGNORE_INDEX = 3
14 TRAIN_IGNORE_INDEX = None
15 NUM_CLASSES = 4
16 MODEL_NAME = "tmp.pytorch"
17 MODEL_SAVE_DIR = "tmp"
18 LOAD = False
19 SAVE = False
20 EPOCHS = 100
21 WEIGHT_DECAY = 0.01
22 TARGET_SIZE = 256
23
24 # Determine the device to use (GPU if available, otherwise CPU)
25 if torch.backends.mps.is_available():
26     device = torch.device("mps")
27 elif torch.cuda.is_available():
28     device = torch.device("cuda")

```



```

29 else:
30     device = torch.device("cpu")
31
32 target_batch_size = 64
33 batch_size = 2
34
35 # Create datasets for training, validation, and testing
36 training_data = dataset("datasets/astrain/color", "datasets/astrain/label", target_transform=target_remap())
37 val_data = dataset("datasets/Val/color", "datasets/Val/label", target_transform=target_remap())
38 test_data = dataset("datasets/Test/color", "datasets/Test/label", target_transform=target_remap())
39
40 # Create data loaders for training, validation, and testing
41 train_dataloader = DataLoader(training_data, batch_size=batch_size, shuffle=True, pin_memory=True)
42 val_dataloader = DataLoader(val_data, batch_size=batch_size, shuffle=True, pin_memory=True, collate_fn=diff_size_collate)
43 test_dataloader = DataLoader(test_data, batch_size=batch_size, shuffle=True, pin_memory=True, collate_fn=diff_size_collate)
44
45
46 # Class Weights
47 class_weight = Tensor([0.30711034803008996, 1.5412496145750956, 1.8445296893647247, 0.30711034803008996])
48 class_weight = Tensor([0.2046795970925636, 1.0271954434416883, 1.2293222812780409, 1.5388026781877073])
49 # class_weight = [1, 1, 1, 1]
50 # class_weight = calculate_class_weights(training_data, 4, None, "dataset")
51 class_weight = class_weight.to(device)
52 class_weight = None
53
54 # Calculate the number of accumulation steps
55 accumulation_steps = target_batch_size // batch_size
56
57 # Model
58 model = unet(3, 4).to(device)
59
60 # Losses
61 train_loss_fn = WeightedDiceCELoss(ignore_index=TRAIN_IGNORE_INDEX, smooth_dice=1, class_weights=class_weight)
62 val_loss_fn = WeightedDiceCELoss(ignore_index=EVAL_IGNORE_INDEX, class_weights=class_weight)
63
64 train_loss_fn = nn.CrossEntropyLoss(weight=class_weight)
65 val_loss_fn = nn.CrossEntropyLoss(weight=class_weight, ignore_index=EVAL_IGNORE_INDEX)
66
67 train_loss_fn = nn.CrossEntropyLoss()
68 val_loss_fn = nn.CrossEntropyLoss(ignore_index=EVAL_IGNORE_INDEX)
69
70 # Optimizer
71 optimizer = optim.AdamW(model.parameters(), weight_decay=WEIGHT_DECAY)
72
73 # Scheduler
74 scheduler = None
75
76 # Metric History
77 agg = MetricsHistory(NUM_CLASSES, EVAL_IGNORE_INDEX)
78
79 # Training Pipeline
80 start(
81     model_save_dir=MODEL_SAVE_DIR,
82     model_save_name=MODEL_NAME,
83     model=model,
84     optimizer=optimizer,
85     train_dataloader=train_dataloader,
86     val_dataloader=val_dataloader,
87     accumulation_steps=accumulation_steps,
88     device=device,
89     train_loss_fn=train_loss_fn,
90     val_loss_fn=val_loss_fn,
91     scheduler=scheduler,
92     agg=agg,
93     load=LOAD,
94     save=SAVE,
95     num_classes=NUM_CLASSES,

```

```

96     ignore_index=EVAL_IGNORE_INDEX,
97     target_size=TARGET_SIZE
98 )

```

Autoencoder training (autoencoder.ipynb)

```

1  from pickle import TRUE
2  from tqdm import tqdm
3  import torch
4  import os
5  import numpy as np
6  from torch.utils.data import DataLoader
7  from torch import nn
8
9  from utils.dataset import *
10 from utils.utils import *
11 from utils.training import trainReconstruction, evalReconstruction
12 from autoencoder.autoencoder import ReconstructionAutoencoder
13
14
15 batch_size = 2
16 target_batch_size = 64
17 accumulation_steps = target_batch_size // batch_size
18
19 training_data = dataset("datasets/astrain/color", "datasets/astrain/label", target_transform=target_remap())
20 validation_data = dataset("datasets/Val/color", "datasets/Val/label", target_transform=target_remap())
21 test_data = dataset("datasets/Test/color", "datasets/Test/label", target_transform=target_remap())
22
23 train_dataloader = DataLoader(training_data, batch_size=batch_size, shuffle=True, pin_memory=True)
24 val_dataloader = DataLoader(validation_data, batch_size=batch_size, shuffle=True, pin_memory=True, collate_fn=diff_size_collate)
25 test_dataloader = DataLoader(test_data, batch_size=batch_size, shuffle=True, pin_memory=True, collate_fn=diff_size_collate)
26
27
28 if torch.backends.mps.is_available():
29     device = torch.device("mps")
30 elif torch.cuda.is_available():
31     device = torch.device("cuda")
32 else:
33     device = torch.device("cpu")
34
35
36 model = ReconstructionAutoencoder(din=3, dout=3).to(device)
37 loss_fn = nn.MSELoss()
38 learning_rate = 1e-3
39 optimizer = torch.optim.Adam(model.parameters(), lr=learning_rate)
40 MODEL_SAVE_DIR = "tmp"
41 MODEL_NAME = "tmp.pytorch"
42 start_epoch = 0
43 recoverCheckpoint = False
44 TARGET_SIZE = 256
45
46 if recoverCheckpoint and os.path.isfile(f"{MODEL_SAVE_DIR}/{MODEL_NAME}"):
47     print(f"Loading checkpoint from: {MODEL_SAVE_DIR}/{MODEL_NAME}")
48     # Load the checkpoint dictionary; move tensors to the correct device
49     checkpoint = torch.load(f"{MODEL_SAVE_DIR}/{MODEL_NAME}", map_location=device, weights_only=False)
50
51     # Load model state
52     model.load_state_dict(checkpoint["model_state_dict"])
53     print(" -> Model state loaded.")
54
55     # Load optimizer state
56     try:
57         optimizer.load_state_dict(checkpoint["optimizer_state_dict"])
58         print(" -> Optimizer state loaded.")
59     except Exception as e:

```

```

60         print(f" -> Warning: Could not load optimizer state: {e}. Optimizer will start from scratch.")
61
62         # Load training metadata
63         start_epoch = checkpoint.get("epoch", 0) # Load last completed epoch, training continues from next one
64         best_val_loss = checkpoint.get("best_val_loss", np.inf)
65
66         print(f" -> Resuming training from epoch {start_epoch + 1}")
67         loaded_notes = checkpoint.get("notes", "N/A")
68         print(f" -> Notes from checkpoint: {loaded_notes}")
69
70     else:
71         print(f"Checkpoint file not found at {MODEL_SAVE_DIR}/{MODEL_NAME}. Starting training from scratch.")
72
73
74
75     best_val_loss = np.inf
76     EPOCHS = 100
77     print("\nStarting Training (Autoencoder)...")
78     for t in range(start_epoch, EPOCHS):
79         current_epoch = t + 1
80         print(f"Epoch {t+1}\n-----")
81         train_loss = trainReconstruction(train_dataloader, model, loss_fn, optimizer, accumulation_steps)
82
83         wrong_val_loss, correct_val_loss = evalReconstruction(val_dataloader, model, loss_fn, target_size=TARGET_SIZE)
84
85         # Save model based on validation val loss improvement
86         if correct_val_loss < best_val_loss:
87             print(f"Validation loss improved ({best_val_loss:.6f} -> {correct_val_loss:.6f}). Saving model...")
88             best_val_loss = correct_val_loss # Save corresponding loss
89             checkpoint_path = os.path.join(MODEL_SAVE_DIR, f"{MODEL_NAME}") # Changed name
90             checkpoint = {
91                 "epoch": t + 1,
92                 "model_state_dict": model.state_dict(),
93                 "optimizer_state_dict": optimizer.state_dict(),
94                 "best_val_loss": best_val_loss,
95             }
96             torch.save(checkpoint, checkpoint_path)
97
98         else:
99             print(f"Corresponding validation loss: {correct_val_loss:.6f} not better than {best_val_loss}")
100
101         #print(f"Wrong Validation loss: {wrong_val_loss:.6f}")
102         print(f"Train loss: {train_loss:.6f}")
103
104         # PLOT a training image reconstruction
105         img, label = training_data[0]
106         img = img.to(device)
107         res = model(img.unsqueeze(0))
108         plt.imshow(res[0].permute(1,2,0).cpu().detach().numpy())
109         plt.savefig(f"drive/MyDrive/autoencoder/images/test{t}.png", format="png")
110         plt.show()
111
112     print("\n--- Training Finished! ---")
113     print(f"Best model saved to: {os.path.join(MODEL_SAVE_DIR, f'{MODEL_NAME}')}")
114
115
116     import torch
117     from torch import Tensor
118     from torch import optim
119     from torch.utils.data import DataLoader
120     from utils.training import start
121     from utils.MetricsHistory import MetricsHistory
122     from utils.weighted_loss import WeightedDiceCELoss
123     from utils.utils import calculate_class_weights
124     from utils.dataset import dataset, target_remap, diff_size_collate
125     from autoencoder.autoencoder import SegmentationAutoencoder
126

```

```

127 EVAL_IGNORE_INDEX = 3
128 TRAIN_IGNORE_INDEX = None
129 NUM_CLASSES = 4
130 MODEL_NAME = "tmp.pytorch"
131 MODEL_SAVE_DIR = "tmp"
132 LOAD = False
133 SAVE = False
134 EPOCHS = 100
135 WEIGHT_DECAY = 0.01
136 TARGET_SIZE = 256
137
138 if torch.backends.mps.is_available():
139     device = torch.device("mps")
140 elif torch.cuda.is_available():
141     device = torch.device("cuda")
142 else:
143     device = torch.device("cpu")
144
145 target_batch_size = 64
146 batch_size = 64
147
148 # With Augmentation
149 training_data = dataset("datasets/astrain/color", "datasets/astrain/label", target_transform=target_remap())
150 val_data = dataset("datasets/Val/color", "datasets/Val/label", target_transform=target_remap())
151 test_data = dataset("datasets/Test/color", "datasets/Test/label", target_transform=target_remap())
152
153 train_dataloader = DataLoader(training_data, batch_size=batch_size, shuffle=True, pin_memory=True)
154 val_dataloader = DataLoader(val_data, batch_size=batch_size, shuffle=True, pin_memory=True, collate_fn=diff_size_collate)
155 test_dataloader = DataLoader(test_data, batch_size=batch_size, shuffle=True, pin_memory=True, collate_fn=diff_size_collate)
156
157
158 # Class Weights
159 class_weight = Tensor([0.33265044664009075, 1.669423957743164, 1.9979255956167454, 0.0])
160 class_weight = Tensor([0.30711034803008996, 1.5412496145750956, 1.8445296893647247, 0.30711034803008996])
161 class_weight = Tensor([0.2046795970925636, 1.0271954434416883, 1.2293222812780409, 1.5388026781877073])
162 # class_weight = [1, 1, 1, 1]
163 # class_weight = calculate_class_weights(training_data, 4, None, "dataset")
164 class_weight = class_weight.to(device)
165
166 accumulation_steps = target_batch_size // batch_size
167
168 # Model
169 # pretrained_encoder_path = "/content/drive/MyDrive/autoencoder/256_with_aug_LR1e-3/checkpoint_256_with_aug_TargetSize256.pth"
170 # model = SegmentationAutoencoder(3, 4, pretrained_encoder_path).to(device)
171 model = SegmentationAutoencoder(3, 4).to(device)
172
173 # Losses
174 train_loss_fn = WeightedDiceCELoss(ignore_index=TRAIN_IGNORE_INDEX, smooth_dice=1, class_weights=class_weight)
175 val_loss_fn = WeightedDiceCELoss(ignore_index=EVAL_IGNORE_INDEX, class_weights=class_weight)
176
177 # train_loss_fn = nn.CrossEntropyLoss(weight=class_weight)
178 # val_loss_fn = nn.CrossEntropyLoss(weight=class_weight, ignore_index=EVAL_IGNORE_INDEX)
179
180 # train_loss_fn = nn.CrossEntropyLoss()
181 # val_loss_fn = nn.CrossEntropyLoss(ignore_index=EVAL_IGNORE_INDEX)
182
183 # Optimizer
184 optimizer = optim.AdamW(model.parameters(), weight_decay=WEIGHT_DECAY)
185
186 # Scheduler
187 scheduler = None
188
189 # Metric History
190 agg = MetricsHistory(NUM_CLASSES, EVAL_IGNORE_INDEX)
191
192 # Training Pipeline
193 start(

```

```

194     model_save_dir=MODEL_SAVE_DIR,
195     model_save_name=MODEL_NAME,
196     model=model,
197     optimizer=optimizer,
198     train_dataloader=train_dataloader,
199     val_dataloader=val_dataloader,
200     accumulation_steps=accumulation_steps,
201     device=device,
202     train_loss_fn=train_loss_fn,
203     val_loss_fn=val_loss_fn,
204     scheduler=scheduler,
205     agg=agg,
206     load=LOAD,
207     save=SAVE,
208     num_classes=NUM_CLASSES,
209     ignore_index=EVAL_IGNORE_INDEX,
210     target_size=TARGET_SIZE
211 )

```

Clip training (clip.ipynb)

```

1  import torch
2  from torch import Tensor
3  from torch import optim
4  from torch.utils.data import DataLoader
5  from utils.training import start
6  from utils.MetricsHistory import MetricsHistory
7  from clip.clipunet import ClipUNet
8  from clip.clipunet_noskips import ClipUNetNoSkips
9  from utils.weighted_loss import WeightedDiceCELoss
10 from utils.utils import calculate_class_weights
11 from utils.dataset import dataset, target_remap, diff_size_collate
12
13 EVAL_IGNORE_INDEX = 3
14 TRAIN_IGNORE_INDEX = None
15 NUM_CLASSES = 4
16 MODEL_NAME = "tmp.pytorch"
17 MODEL_SAVE_DIR = "tmp"
18 LOAD = False
19 SAVE = False
20 EPOCHS = 100
21 WEIGHT_DECAY = 0.01
22 PRETRAINED_MODEL_NAME = "openai/clip-vit-base-patch16"
23 TARGET_SIZE = 224
24 SKIP_LAYER_INDICES = [3, 5, 7, 9]
25
26 if torch.backends.mps.is_available():
27     device = torch.device("mps")
28 elif torch.cuda.is_available():
29     device = torch.device("cuda")
30 else:
31     device = torch.device("cpu")
32
33 # Class Weights
34 class_weight = Tensor([0.33265044664009075, 1.669423957743164, 1.9979255956167454, 0.0])
35 class_weight = Tensor([0.30711034803008996, 1.5412496145750956, 1.8445296893647247, 0.30711034803008996])
36 class_weight = Tensor([0.2046795970925636, 1.0271954434416883, 1.2293222812780409, 1.5388026781877073])
37 # class_weight = [1, 1, 1, 1]
38 # class_weight = calculate_class_weights_v3(training_data, 4, None, "dataset")
39 class_weight = class_weight.to(device)
40
41 target_batch_size = 64
42 batch_size = 2
43
44 training_data = dataset("datasets/rstrain/color", "datasets/rstrain/label", target_transform=target_remap())

```

```

45 val_data = dataset("datasets/Val/color", "datasets/Val/label", target_transform=target_remap())
46 test_data = dataset("datasets/Test/color", "datasets/Test/label", target_transform=target_remap())
47
48 train_dataloader = DataLoader(training_data, batch_size=batch_size, shuffle=True, pin_memory=True)
49 val_dataloader = DataLoader(val_data, batch_size=batch_size, shuffle=True, pin_memory=True, collate_fn=diff_size_collate)
50 test_dataloader = DataLoader(test_data, batch_size=batch_size, shuffle=True, pin_memory=True, collate_fn=diff_size_collate)
51
52 accumulation_steps = target_batch_size // batch_size
53
54 # Model
55 model = ClipUNet().to(device)
56
57 # Losses
58 train_loss_fn = WeightedDiceCELoss(ignore_index=TRAIN_IGNORE_INDEX, smooth_dice=1, class_weights=class_weight)
59 val_loss_fn = WeightedDiceCELoss(ignore_index=EVAL_IGNORE_INDEX, class_weights=class_weight)
60
61 # Optimizer
62 optimizer = optim.AdamW(model.parameters(), weight_decay=WEIGHT_DECAY)
63
64 # Scheduler
65 scheduler = None
66
67 # Metric History
68 agg = MetricsHistory(NUM_CLASSES, EVAL_IGNORE_INDEX)
69
70 # Training Pipeline
71 start(
72     model_save_dir=MODEL_SAVE_DIR,
73     model_save_name=MODEL_NAME,
74     model=model,
75     optimizer=optimizer,
76     train_dataloader=train_dataloader,
77     val_dataloader=val_dataloader,
78     accumulation_steps=accumulation_steps,
79     device=device,
80     train_loss_fn=train_loss_fn,
81     val_loss_fn=val_loss_fn,
82     scheduler=scheduler,
83     agg=agg,
84     load=LOAD,
85     save=SAVE,
86     num_classes=NUM_CLASSES,
87     ignore_index=EVAL_IGNORE_INDEX,
88     target_size=TARGET_SIZE
89 )

```

Prompt based model training (prompt.ipynb)

```

1 import torch
2 from torch import Tensor
3 from torch import optim
4 from torch.utils.data import DataLoader
5 from prompt_based.prompt import PromptModel
6 from utils.MetricsHistory import MetricsHistory
7 from utils.weighted_loss import WeightedDiceNLLLoss
8 from utils.utils import calculate_class_weights
9 from utils.dataset import promptDataset, diff_size_collate
10 from utils.training import start_prompt
11
12 EVAL_IGNORE_INDEX = 3
13 TRAIN_IGNORE_INDEX = None
14 NUM_CLASSES = 4
15 MODEL_NAME = "tmp.pytorch"
16 MODEL_SAVE_DIR = "tmp"
17 LOAD = False

```

```

18 SAVE = False
19 EPOCHS = 100
20 WEIGHT_DECAY = 0.01
21 PRETRAINED_MODEL_NAME = "openai/clip-vit-base-patch16"
22 TARGET_SIZE = 224
23 SKIP_LAYER_INDICES = [3, 5, 7, 9]
24 CLIP_PATH="/content/drive/MyDrive/clip/runs/clip_256_ce_dice_full_weight_fix_train_eval.pytorch"
25
26 if torch.backends.mps.is_available():
27     device = torch.device("mps")
28 elif torch.cuda.is_available():
29     device = torch.device("cuda")
30 else:
31     device = torch.device("cpu")
32
33 # Class Weights
34 class_weight = Tensor([0.33265044664009075, 1.669423957743164, 1.9979255956167454, 0.0])
35 class_weight = Tensor([0.30711034803008996, 1.5412496145750956, 1.8445296893647247, 0.30711034803008996])
36 class_weight = Tensor([0.2046795970925636, 1.0271954434416883, 1.2293222812780409, 1.5388026781877073])
37 class_weight = Tensor([1, 1, 1, 1])
38 # class_weight = calculate_class_weights_v3(training_data, 4, None, "dataset")
39 class_weight = class_weight.to(device)
40
41 target_batch_size = 64
42 batch_size = 2
43
44 training_data = promptDataset("datasets/pstrain/color", "datasets/pstrain/point_prompt", "datasets/pstrain/label")
45 val_data = promptDataset("datasets/psVal/color", "datasets/psVal/point_prompt", "datasets/psVal/label")
46 test_data = promptDataset("datasets/psTest/color", "datasets/psTest/point_prompt", "datasets/psTest/label")
47
48 train_dataloader = DataLoader(training_data, batch_size=batch_size, shuffle=True, pin_memory=True)
49 val_dataloader = DataLoader(val_data, batch_size=batch_size, shuffle=True, pin_memory=True, collate_fn=diff_size_collate)
50 test_dataloader = DataLoader(test_data, batch_size=batch_size, shuffle=True, pin_memory=True, collate_fn=diff_size_collate)
51
52 accumulation_steps = target_batch_size // batch_size
53
54 # Model
55 # model = PromptModel(CLIP_PATH).to(device) # oad pretrained clip
56 model = PromptModel().to(device)
57
58 # Losses
59 stable_log = lambda x: torch.log(x + 1e-9)
60 train_loss_fn = WeightedDiceNLLLoss(ignore_index=TRAIN_IGNORE_INDEX, smooth_dice=1, class_weights=class_weight, apply_softmax=True)
61 val_loss_fn = WeightedDiceNLLLoss(ignore_index=EVAL_IGNORE_INDEX, class_weights=class_weight, apply_softmax=False, nll_loss=stable_log)
62
63 # Optimizer
64 optimizer = optim.AdamW(model.parameters(), weight_decay=WEIGHT_DECAY)
65
66 # Scheduler
67 scheduler = None
68
69 # Metric History
70 agg = MetricsHistory(NUM_CLASSES, EVAL_IGNORE_INDEX)
71
72 # Training Pipeline
73 start_prompt(
74     model_save_dir=MODEL_SAVE_DIR,
75     model_save_name=MODEL_NAME,
76     model=model,
77     optimizer=optimizer,
78     train_dataloader=train_dataloader,
79     val_dataloader=val_dataloader,
80     accumulation_steps=accumulation_steps,
81     device=device,
82     train_loss_fn=train_loss_fn,
83     val_loss_fn=val_loss_fn,
84     scheduler=scheduler,

```

```
85     agg=agg,  
86     load=LOAD,  
87     save=SAVE,  
88     num_classes=NUM_CLASSES,  
89     ignore_index=EVAL_IGNORE_INDEX,  
90     target_size=TARGET_SIZE  
91 )
```
