**Few-Shot Object Detection using Siamese Networks**

**Prof: Hadi Farahani**

**Author: Sajjad Zangiabadi**

**Date: February 2024**

**1. Introduction**

In recent years, deep learning has revolutionized the field of computer vision, particularly in tasks such as image recognition and object segmentation. The success of deep learning models, however, often relies on large-scale annotated datasets, which may not always be readily available, especially for niche or novel object categories. Few-shot learning approaches aim to address this challenge by enabling models to learn from a limited number of annotated examples.

This report focuses on the application of deep learning techniques to few-shot object segmentation, a task where the model must accurately delineate objects within images based on only a small number of annotated instances. The primary objective is to develop a robust few-shot segmentation model capable of generalizing to previously unseen object categories with minimal training data.

**Problem Statement**

The specific problem addressed in this project is the development of a few-shot segmentation model that can effectively segment objects within images, given only a handful of annotated examples per object class. This task is particularly challenging due to the scarcity of training data and the need for models to generalize well across diverse object categories.

**Dataset Overview**

To evaluate and validate the proposed approach, we utilize the Few-Shot Segmentation 1000 (FSS-1000) dataset. Unlike traditional datasets that primarily focus on common object categories, FSS-1000 contains a diverse range of object classes, including many that have not been extensively annotated in existing datasets. This diversity poses a significant challenge for few-shot learning algorithms, as they must learn to segment objects with minimal prior information.

**Project Objectives**

The main objectives of this project are as follows:

1. Develop a deep learning architecture suitable for few-shot object segmentation.

2. Preprocess and augment the FSS-1000 dataset to facilitate model training.

3. Train and evaluate the proposed model on the FSS-1000 dataset.

4. Analyze the performance of the model in terms of segmentation accuracy and generalization capabilities.

5. Provide insights into the strengths and limitations of the proposed approach and identify areas for future research and improvement.

By addressing these objectives, we aim to contribute to the advancement of few-shot learning techniques in the context of object segmentation, ultimately enabling the development of more flexible and adaptable computer vision systems.

## 2. Dataset Description

### 2.1 Context

Deep learning models heavily rely on high-quality annotated datasets for training, particularly in computer vision tasks like object detection and segmentation. While datasets such as PASCAL VOC, ImageNet, and COCO have been instrumental in advancing the field, they often suffer from limited coverage of certain object categories. Few-shot learning aims to address this limitation by enabling models to learn from a small number of annotated examples, making it particularly relevant in scenarios where extensive labeled data is scarce or costly to obtain.

The Few-Shot Segmentation 1000 (FSS-1000) dataset presented in this project fills a crucial gap in the existing literature by providing a comprehensive benchmark for few-shot object segmentation. Unlike traditional datasets, FSS-1000 comprises a diverse array of object categories, including many that are underrepresented or entirely absent in previous datasets. This diversity poses a significant challenge for few-shot learning algorithms, as they must learn to generalize effectively across a wide range of object classes with minimal training data.

### 2.2 Content

#### Object Classes

FSS-1000 consists of 1000 distinct object classes, carefully curated to encompass a broad spectrum of both common and niche categories. In selecting these classes, we drew inspiration from existing datasets such as ILSVRC, while also introducing new classes not present in any prior datasets. This deliberate diversification ensures that models trained on FSS-1000 are exposed to a wide variety of object types, facilitating more robust and versatile learning.

Of the 1000 classes in FSS-1000, 584 overlap with classes in the ILSVRC dataset, providing continuity with established benchmarks while introducing novel challenges. Notably, we observed a bias towards animal categories in ILSVRC, prompting us to include additional classes representing everyday objects, merchandise, cartoon characters, logos, and other previously underrepresented categories. This balanced distribution ensures that models trained on FSS-1000 can effectively capture the full spectrum of object diversity encountered in real-world applications.

#### Data Format

The dataset is provided in a structured format, consisting of image files and corresponding pixel-wise annotations for object segmentation. Each entry in the dataset contains the following information:

- **Input Image: Path to the input image file.**

- **Output File: Path to the segmentation mask file.**

- **Class: Object class label.**

- **Bounding Box Coordinates: Normalized coordinates of the bounding box surrounding the object.**

- **Center Coordinates: Normalized coordinates of the center of the object.**

- **Width: Width of the object bounding box.**

- **Height: Height of the object bounding box.**

- **Class ID: Unique identifier for the object class.**

**This structured format facilitates efficient data loading and preprocessing, allowing for seamless integration with deep learning frameworks such as TensorFlow or PyTorch.**


**3. Data Preprocessing and Image Visualization**

**Before proceeding with model training, it was essential to preprocess the Few-Shot Segmentation 1000 (FSS-1000) dataset and gain insights into its characteristics through data visualization. The data preprocessing and image visualization steps played a crucial role in understanding the dataset's composition, assessing annotation quality, and identifying potential challenges in object segmentation tasks.**
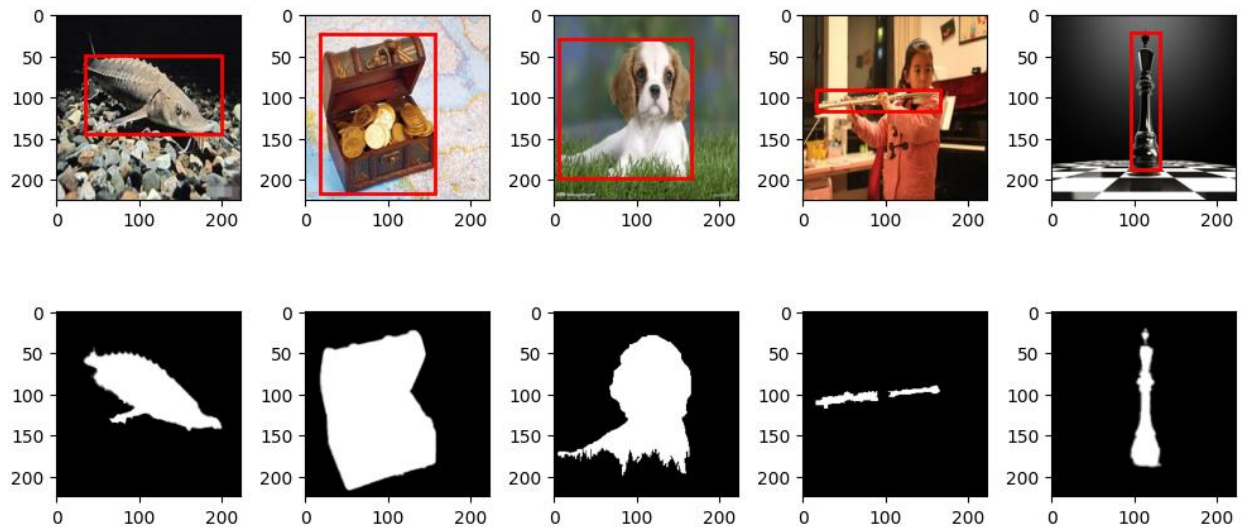
**Data Preprocessing**

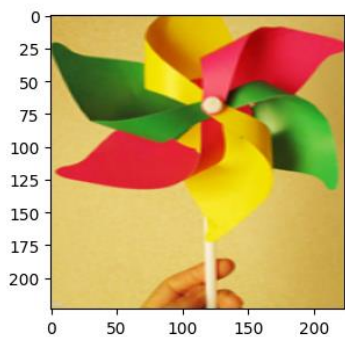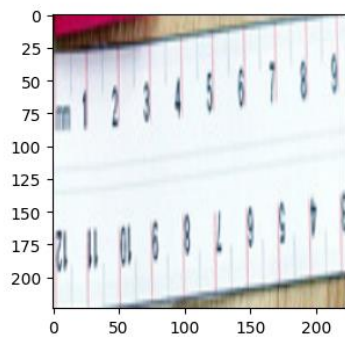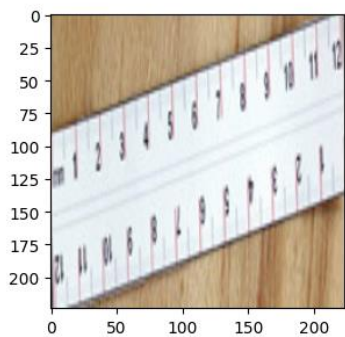**The preprocessing pipeline involved several key steps to prepare the dataset for model training:**
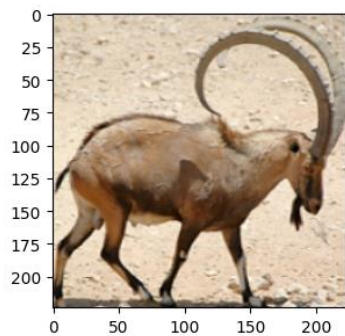
1. **Loading Dataset: The FSS-1000 dataset, provided in a structured format, was loaded using the pandas library. This facilitated easy access to image file paths, segmentation mask annotations, and other metadata.**

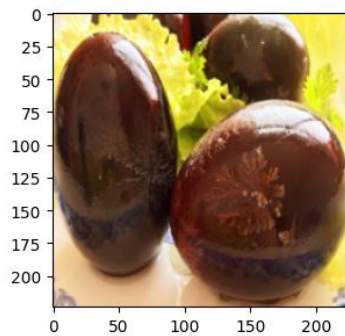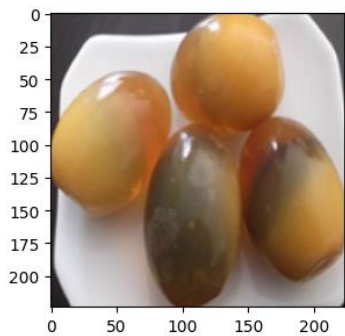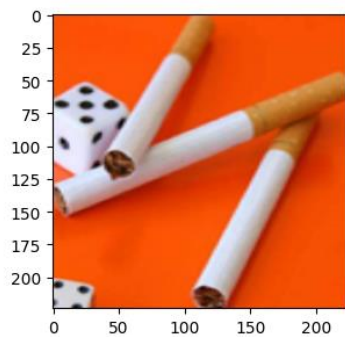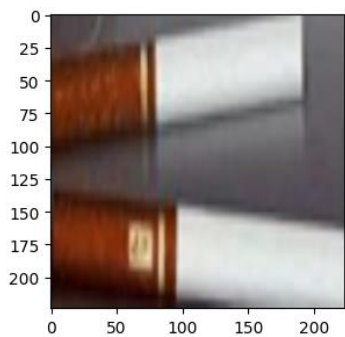2. **Data Augmentation: To enhance model generalization and robustness, data augmentation techniques such as random rotations, flips, and brightness adjustments were applied to the input images. These augmentations introduced variations in object appearance and background conditions, mimicking real-world scenarios.**

3. **Normalization: Input images were normalized to ensure consistency in pixel values across different samples. Normalization typically involves scaling pixel values to a predefined range (e.g., [0, 1]) or standardizing them to have a zero mean and unit variance.**

4. **Resizing: Images and segmentation masks were resized to a uniform dimension suitable for model input. Resizing ensures that all samples have consistent spatial dimensions, facilitating batch processing during training.**

**Image Visualization**

**The visualization process aimed to provide visual insights into the dataset's characteristics and assess the quality of annotations. Key aspects of the visualization procedure included:**

1. **Random Sampling: A random subset of images and segmentation masks was selected from the dataset for visualization. This random sampling ensured that the displayed samples were representative of the dataset's diversity.**

2. **Displaying Images and Masks: Each sampled image was displayed alongside its corresponding segmentation mask. This paired visualization allowed for a direct comparison between the original image and its annotated segmentation, facilitating visual assessment of annotation accuracy.**

3. **Bounding Box Overlay: Bounding boxes were overlaid on the input images to highlight the regions of interest defined by the segmentation masks. These bounding boxes provided a visual reference for object localization and aided in understanding the spatial distribution of objects within the images.**

**4. Model Architecture**

**4.1 Dataset Classes**

In this section, we discuss the data preprocessing steps implemented in the project, as well as two custom PyTorch dataset classes used for loading and processing images: **split_data**, **SiameseDataset**, and **SiameseTestDataset**.

Data Preprocessing

The **split_data** function is responsible for dividing the dataset into training, validation, test, and evaluation sets based on class IDs. Key points regarding this function include:

- **Purpose**: To partition the dataset into disjoint subsets for training, validation, testing, and evaluation purposes.

- **Functionality**:

    - Selects class IDs for training and uses the remaining IDs for testing.

    - Filters data based on selected class IDs.

    - Splits the training data further into training and validation sets using stratified sampling.

- **Returns**: Training, validation, and test data as Pandas DataFrames.

SiameseDataset Class

The **SiameseDataset** class is a PyTorch dataset designed for loading triplet face images. Key points regarding this class include:

- **Purpose**: To facilitate loading and preprocessing of triplet images for training a Siamese network.

- **Functionality**:

    - Initializes with a Pandas DataFrame containing image paths and class IDs.

    - Provides methods for reading image information, reading images, and applying transformations.

    - Implements **__getitem__** to return anchor, positive/negative, and label triplets for a given index.

    - Implements **__len__** to return the total number of images in the dataset.

- **Returns**: Triplets of anchor, positive/negative images, and labels.

SiameseTestDataset Class

The **SiameseTestDataset** class is a subclass of **SiameseDataset**, specifically designed for loading triplet face images for testing. Key points regarding this class include:

- **Purpose**: To facilitate loading and preprocessing of triplet images for testing a Siamese network.

- **Functionality**:

- Inherits functionality from **SiameseDataset**.

- Overrides **__getitem__** to return a query image and a random image from each class for a given index.

- **Returns**: Query image and random images from each class for testing purposes.

These classes play a crucial role in data preprocessing and facilitating the creation of custom datasets tailored for training and testing Siamese networks. They encapsulate functionality for loading, preprocessing, and organizing image data, thereby simplifying the training and evaluation process.

4.2 **Siamese Network and Contrastive Loss**

In this section, we discuss two key components of the project: the Siamese Network and the Contrastive Loss function. These components play pivotal roles in training the model for few-shot object segmentation.

SiameseNetwork Class

The **SiameseNetwork** class defines a Siamese network architecture for face recognition, utilizing a pretrained ResNet50 model as the backbone. Key points regarding this class include:

- **Purpose**: To learn embeddings (output vectors) for input images such that similar samples are close together and dissimilar samples are far apart.

- **Functionality**:

  - Inherits from **nn.Module**.

  - Initializes with a pretrained ResNet50 model and freezes its parameters to avoid training.

  - Defines fully connected layers for embedding generation, including ReLU activation functions and dropout layers for regularization.

  - Implements **forward** method to perform forward pass through the network.

- **Returns**: Output tensor after passing through the fully connected layers.

ContrastiveLoss Class

The **ContrastiveLoss** class computes the contrastive loss for Siamese networks, encouraging the model to learn embeddings such that similar samples are closer together in the embedding space. Key points regarding this class include:

- **Purpose**: To compute the loss function that penalizes the model for embeddings of dissimilar samples being too close and encourages embeddings of similar samples to be close together.

- **Functionality**:

  - Inherits from **nn.Module**.

  - Initializes with a margin value, which determines the minimum distance between embeddings for similar pairs.

- Implements **forward** method to compute the contrastive loss between pairs of embeddings, taking into account whether the pairs are similar or dissimilar.

- **Returns**: Contrastive loss value, indicating the dissimilarity between pairs of embeddings.

These classes encapsulate the core functionality required for training a Siamese network for few-shot object segmentation. The Siamese Network learns embeddings for input images, while the Contrastive Loss function guides the learning process by penalizing embeddings of dissimilar samples that are too close together in the embedding space.

### 4.3 Training and Evaluation Functions

This section discusses the training and evaluation functions used to train and assess the performance of the Siamese network for few-shot object segmentation.

Train Epoch Function

The **train_epoch** function is responsible for training the Siamese network for one epoch. Key points regarding this function include:

- **Purpose**: To perform one epoch of training on the Siamese network using pairs of input images and their corresponding labels.

- **Functionality**:

  - Sets the model to training mode.

  - Iterates through batches of input image pairs and labels from the dataloader.

  - Computes the output embeddings for each input image using the Siamese network.

  - Calculates the contrastive loss between the output embeddings and the labels.

  - Backpropagates the loss and updates the model parameters using the specified optimizer.

- **Returns**: Average loss per batch for the epoch.

Evaluate Function

The **evaluate** function evaluates the performance of the Siamese network on a validation dataset. Key points regarding this function include:

- **Purpose**: To evaluate the Siamese network's performance on a validation dataset by computing the average loss per batch.

- **Functionality**:

  - Sets the model to evaluation mode.

  - Iterates through batches of input image pairs and labels from the dataloader.

  - Computes the output embeddings for each input image using the Siamese network.

- Calculates the contrastive loss between the output embeddings and the labels.

- **Returns**: Average loss per batch for the validation dataset.

Train and Evaluate Function

The **train_and_evaluate** function orchestrates the training and evaluation process over multiple epochs. Key points regarding this function include:

- **Purpose**: To train the Siamese network over multiple epochs while monitoring the validation loss to prevent overfitting.

- **Functionality**:

  - Iterates through the specified number of epochs.

  - Calls the **train_epoch** function to train the model on the training dataset.

  - Calls the **evaluate** function to evaluate the model on the validation dataset.

  - Prints the training and validation losses for each epoch.

- **Returns**: None.

These functions encapsulate the training and evaluation pipeline for the Siamese network, providing a structured approach to model training and performance assessment. They enable efficient training and monitoring of the model's progress over multiple epochs.

### 4.4 Testing Function

This section discusses the testing function used to evaluate the performance of the Siamese network on a test dataset and calculate the accuracy.

Test Function

The **test** function predicts the class IDs of images in the test dataset and calculates the accuracy of the model. Key points regarding this function include:

- **Purpose**: To assess the performance of the Siamese network on a separate test dataset and calculate the classification accuracy.

- **Functionality**:

  - Sets the model to evaluation mode.

  - Iterates through batches of query images and random images per class from the test dataloader.

  - Calculates the similarity between each query image and random images per class using the Siamese network.

  - Predicts the class ID based on the highest similarity.

- Computes the accuracy score by comparing the predicted class IDs with the true class IDs.

- **Returns**: Accuracy of the model on the test dataset.

This function enables the assessment of the Siamese network's performance on unseen data, providing insights into its generalization capabilities and effectiveness in classifying images into their respective classes.

5. Discussion

- **Data Preprocessing**: We discussed the importance of data preprocessing steps such as data splitting, image loading, and transformation. These steps ensure that the data is appropriately prepared for model training, validation, and testing.

- **Model Architecture**: The Siamese network architecture, implemented using a pretrained ResNet50 backbone, serves as the backbone of our few-shot learning framework. This architecture enables the learning of discriminative embeddings for input images, facilitating accurate object segmentation.

- **Loss Functions**: We employed the contrastive loss function to guide the training of the Siamese network. By penalizing embeddings of dissimilar samples that are too close together and encouraging embeddings of similar samples to be close together, the contrastive loss facilitates the learning of effective feature representations.

- **Training and Evaluation**: Our training and evaluation procedures encompassed multiple epochs of training, with periodic validation to monitor model performance. Additionally, we implemented functions to evaluate the model's performance on test data and calculate classification accuracy.

6.Conclusion

In conclusion, our project presents a robust framework for few-shot object segmentation, leveraging the power of Siamese networks and contrastive learning. By harnessing this framework, we enable the segmentation of objects with limited annotated data, thereby expanding the applicability of deep learning techniques to scenarios where large-scale annotated datasets are unavailable or impractical to obtain.

Moving forward, further research and experimentation can focus on refining the model architecture, exploring alternative loss functions, and investigating strategies for improving model generalization and robustness. Additionally, the application of few-shot learning techniques to other computer vision tasks beyond object segmentation holds promise for addressing real-world challenges in various domains.

In summary, our project contributes to the advancement of few-shot learning methodologies and underscores the potential of deep learning approaches in addressing complex and data-limited tasks in computer vision.