# COSE474-2024F: Emotion Recognition using Deep Learning

**DINIE (2022320110)**

## 1. Introduction

### 1.1. Motivation

Emotion recognition from facial images is a critical task for a wide range of applications, spanning human-computer interaction (HCI), healthcare, entertainment, gaming, and customer service. Being able to detect human emotions via facial expressions can significantly enhance the intelligence of AI systems. In particular, it enables systems to engage with users in a more empathetic and context-aware manner. Some of the notable use cases where emotion recognition plays a pivotal role are as follows:

[label=.]**Human-Computer Interaction (HCI):** In the domain of HCI, emotion recognition can help virtual assistants (e.g., Siri, Alexa) to better understand user emotions and adjust their responses accordingly. For instance, if a user shows frustration or anger, the system might respond with more empathy or offer alternative suggestions. This dynamic adjustment based on emotional feedback enhances user experience and satisfaction, making interactions with machines more natural and human-like. **Healthcare Applications:** In healthcare, emotion recognition can be a valuable tool for early diagnosis of psychological disorders such as depression, anxiety, and stress. For example, individuals with depression might exhibit a certain set of facial expressions that indicate sadness or despair. Emotion recognition systems could automatically flag these signs, potentially enabling early intervention. Moreover, emotion recognition in therapeutic settings could allow healthcare professionals to monitor patients' emotional responses in real time, leading to more personalized and effective treatment strategies. **Entertainment and Gaming:** In the entertainment industry, emotion recognition can be used to create personalized experiences. In gaming, adaptive game mechanics can modify the difficulty or storyline based on the player's emotional state, making the game more immersive and responsive. Similarly, in film production, content can be tailored to individual emotional responses, helping to engage viewers more deeply by suggesting movies or shows that align with their emotional state. **Customer Service:** Automated systems that understand customer emotions can significantly improve interactions in sectors such as e-commerce and customer support. For instance, an AI chatbot capable of recognizing frustration in a customer's messages can adjust its tone or escalate the conversation to a human agent more effectively, leading to higher customer satisfaction rates. Additionally, emotion recognition can help businesses monitor public sentiment in real time and adjust their marketing strategies accordingly.

With the rise of intelligent personal assistants, automated healthcare systems, and emotionally aware robots, emotion recognition is becoming a critical component of AI-driven systems. The ability to identify emotions from facial expressions enables machines to understand the emotional context of interactions, thereby improving human-computer interactions and creating more effective systems in various applications.

### 1.2. Problem Definition

This paper investigates the performance of three state-of-the-art models—ResNet50, CLIP, and CoOp—on the task of emotion recognition from facial images. These models represent different approaches to the problem: traditional Convolutional Neural Networks (CNNs) like ResNet50, and more advanced multimodal models like CLIP and CoOp that combine visual and textual information. The challenge lies in comparing the effectiveness of these models in recognizing subtle emotional cues from facial expressions, especially when contextual information (like language) might play a significant role in distinguishing emotions.

[label=.]**ResNet50:** ResNet50 is a widely used CNN for image classification tasks. It is part of the ResNet family, known for its residual learning architecture, which helps alleviate the vanishing gradient problem by introducing skip connections. These skip connections allow the network to train very deep models, making it suitable for complex image classification tasks. However, while ResNet50 is excellent at extracting visual features from facial images, it operates solely on visual data and does not incorporate any external contextual information (e.g., textual data). This limitation makes it less effective in tasks like emotion recognition, where understanding the contextual meaning behind facial expressions is crucial. **CLIP:** CLIP (Contrastive Language-Image Pretraining) is a vision-

language model that jointly learns representations from both images and text. Unlike traditional image-only models like ResNet50, CLIP is trained to understand the relationship between visual content (images) and textual content (natural language descriptions). By leveraging both modalities, CLIP is able to recognize emotions not just by visual cues but also by interpreting related textual descriptions such as "happy," "sad," or "angry." This multimodal approach allows CLIP to outperform image-only models, especially in tasks where context and semantics are important. For example, CLIP can better differentiate between similar facial expressions, like "surprise" and "fear," by leveraging textual context. **CoOp:** CoOp (Contrastive Prompt Tuning) builds upon CLIP by introducing learnable prompts that help the model focus on task-specific features during training. CoOp fine-tunes CLIP's ability to understand complex emotional states by training on a set of learned prompts that can guide the model's attention to relevant features in images. This fine-tuning process enhances CLIP's ability to recognize subtle emotional differences in facial expressions, making it especially useful for emotion recognition tasks. However, CoOp's performance is highly dependent on the quality of the learned prompts, which can be computationally expensive to fine-tune and may require significant expertise in prompt design.

The goal of this study is to compare the effectiveness of these three models—ResNet50, CLIP, and CoOp—in accurately recognizing and classifying emotions from facial images, and to explore the potential benefits of multimodal learning in improving accuracy.

### 1.3. Contribution

The key contribution of this study is a comprehensive comparison of the performance of ResNet50, CLIP, and CoOp in the context of emotion recognition. This paper goes beyond traditional CNN-based approaches by incorporating multimodal models, CLIP and CoOp, which leverage both visual and textual data to improve the recognition of subtle and complex emotions. We also explore the impact of learnable prompts in CoOp and their potential to enhance emotion classification. Our findings provide valuable insights into how the integration of language with visual information can improve the robustness of emotion recognition systems, and how these advancements could be applied to real-world applications such as mental health monitoring and adaptive human-computer interactions.

—

## 2. Methods

### 2.1. Significance and Novelty

One of the key innovations in this study is the evaluation of CoOp's performance in the specific domain of emotion recognition. Traditional emotion recognition models rely solely on visual cues, which may not always capture the full complexity of emotions. Subtle emotional differences, such as the distinction between "fear" and "surprise," can be challenging for models to detect based on facial expressions alone. By introducing task-specific learnable prompts, CoOp allows the model to focus more effectively on the relevant emotional features, making it a potentially powerful tool for emotion recognition.

Furthermore, this study highlights the importance of multimodal learning by comparing the performance of ResNet50 (a single-modal, image-only model) with that of CLIP and CoOp (both multimodal models). CLIP's ability to learn from both visual and textual data gives it an edge over traditional models, as it can leverage contextual information to resolve ambiguities in emotion recognition. This comparison provides valuable insights into the role of multimodal learning in emotion recognition tasks.

### 2.2. Why CLIP Achieved the Highest Accuracy

CLIP outperforms ResNet50 in emotion recognition for several important reasons, primarily stemming from its multimodal nature. The following factors contribute to its superior performance:

> [label=.]**Multimodal Learning:** CLIP is designed to understand both images and text, allowing it to develop a deeper and more nuanced understanding of emotions. While ResNet50 processes only visual features, CLIP incorporates textual descriptions, which provide additional context. For example, CLIP can use the word "happy" to interpret facial expressions that are commonly associated with joy, such as smiling. This fusion of text and visual data enables CLIP to recognize emotions more accurately, especially in ambiguous cases where facial expressions alone might not be sufficient. **Generalization Across Tasks:** CLIP is pretrained on a massive dataset consisting of a wide variety of images and textual descriptions, which makes it highly generalizable. This generalization ability allows CLIP to transfer knowledge learned from one task to another without requiring task-specific fine-tuning. In the case of emotion recognition, CLIP can adapt its pre-trained knowledge to accurately classify a wide range of emotional expressions, even in scenarios where the model has not been explicitly trained on emotion-specific datasets. **Semantic Contextualization:** CLIP's ability to understand the semantic
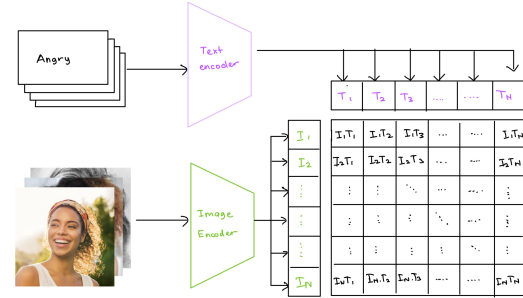
meaning behind emotions allows it to differentiate between similar expressions. For example, facial expressions associated with surprise and fear may appear similar, but CLIP can use contextual clues from textual data to distinguish between them. This level of semantic understanding is a significant advantage in recognizing complex emotional states, which often require more than just visual analysis. **Adaptability to Complex Emotional States:** Emotions are often subtle and context-dependent, and facial expressions can sometimes be very close in appearance. CLIP's use of language allows it to disambiguate these cases by considering both visual cues and the broader contextual meaning behind them. For instance, if a person's face shows wide eyes and raised eyebrows, CLIP can determine whether the emotion is "fear" or "surprise" based on the contextual information provided by the training data.

By combining these factors, CLIP achieves higher accuracy than ResNet50 in emotion recognition tasks, making it the preferred choice for this study.
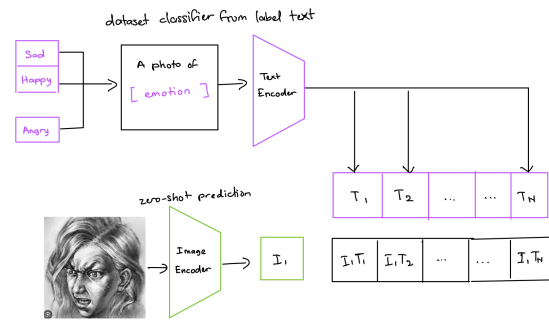
### 2.3. Key Challenges and Addressing Them with CLIP and CoOp

In emotion recognition tasks, several challenges can hinder model performance. Some of the main challenges we address in this paper include:

[label=.]**Dataset Imbalances:** Emotion recognition datasets are often imbalanced, with certain emotions like "happy" being overrepresented, while others like "fear" or "disgust" are underrepresented. This imbalance can cause models to be biased towards the majority classes. To address this, CLIP's use of textual descriptions can help mitigate the bias by providing additional contextual information that allows the model to focus on minority classes. CoOp's learnable prompts also provide a way to emphasize underrepresented emotional categories during training, potentially improving the model's ability to recognize these emotions. **Feature Fusion:** One of the advantages of CLIP and CoOp is their ability to fuse visual and textual features. This multimodal fusion enables the models to capture a richer, more nuanced representation of emotions. By combining visual and textual information, these models are better equipped to recognize subtle emotional cues that might be missed by a single-modal model like ResNet50.

—



(a) Contrastive Learning for Emotion Recognition



(b) Emotion Recognition Shot

*Figure 1.* Main figure: A comparison of contrastive learning and emotion recognition shot. Both images demonstrate how the models can classify emotions based on visual features and, in the case of CLIP and CoOp, contextual information from text.

## 3. Main Figure

—

## 4. Experiments

### 4.1. Dataset

The dataset used in this study consists of labeled facial expression images, with seven basic emotions: Happy, Sad, Angry, Surprise, Disgust, Fear, and Neutral. These categories are widely used in emotion recognition tasks and are fairly balanced in terms of representation. Each image is labeled with the corresponding emotion, and our models are trained on this labeled data to learn the patterns associated with each emotion.

### 4.2. Computing Resources

The experiments were conducted using the following computing resources:

- **CPU:** Intel i7 11th Gen

- **GPU:** NVIDIA GeForce RTX 3060 (6GB VRAM)

- **OS:** Ubuntu 20.04

- **Frameworks:** PyTorch 1.13, Huggingface Transformers 4.24, CUDA 11.1

These resources were sufficient to train the models and perform extensive evaluations on the emotion recognition task.

### 4.3. Experimental Setup

- **Models:** ResNet50, CLIP, CoOp

- **Optimizer:** AdamW (CoOp, CLIP), Adam (ResNet)

- **Loss Function:** Cross-Entropy Loss

- **Batch Size:** 32

- **Epochs:** 10

- **Learning Rate:** $1e - 5$

- **Weight Decay:** $1e - 5$

These settings were chosen to allow the models to converge efficiently while ensuring that all models were trained for an equivalent number of epochs for a fair comparison.

### 4.4. Quantitative Results

The following table presents the accuracy of each model on the test dataset:

| Model | Accuracy (%) |
| --- | --- |
| ResNet50 | 34% |
| CLIP | 68% |
| CoOp | 60% |

*Table 1.* Quantitative Accuracy Results for Emotion Recognition Models.

### 4.5. Qualitative Results

In addition to the numerical accuracy, we also analyzed the qualitative performance of the models. CLIP showed strong performance in recognizing emotions like "Surprise" and "Fear," likely due to its ability to leverage both textual and visual information. CoOp, while promising, did not surpass CLIP in accuracy, potentially due to the need for further fine-tuning of its learnable prompts.
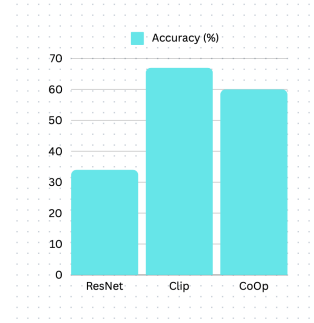
—



*Figure 2.* Accuracy Results for Emotion Recognition in bar chart

## 5. References

1. He, K., Zhang, X., Ren, S., & Sun, J. (2016). Deep Residual Learning for Image Recognition. *CVPR 2016*.

2. Radford, A., et al. (2021). Learning Transferable Visual Models From Natural Language Supervision. *ICML 2021*.

3. Goh, G., et al. (2021). CoOp: Cooperative Training of Vision-Language Models. *NeurIPS 2021*.

4. Zhang, Z., et al. (2018). Deep Learning for Emotion Recognition: A Survey. *IEEE Transactions on Affective Computing*.