# Paper Reading

Learning Transferable Visual Models From Natural Language Supervision
Contrastive Language-Image Pre-training

# Introduction

• The paper addresses the challenge of training computer vision models that require large amounts of labeled data, which is both costly and time-consuming

• It introduces CLIP (Contrastive Language-Image Pre-training), a model that learns visual concepts from image-text pairs using natural language supervision

• The key idea is to train a model that learns visual representations directly from text, enabling zero-shot transfer to various downstream tasks without task-specific training
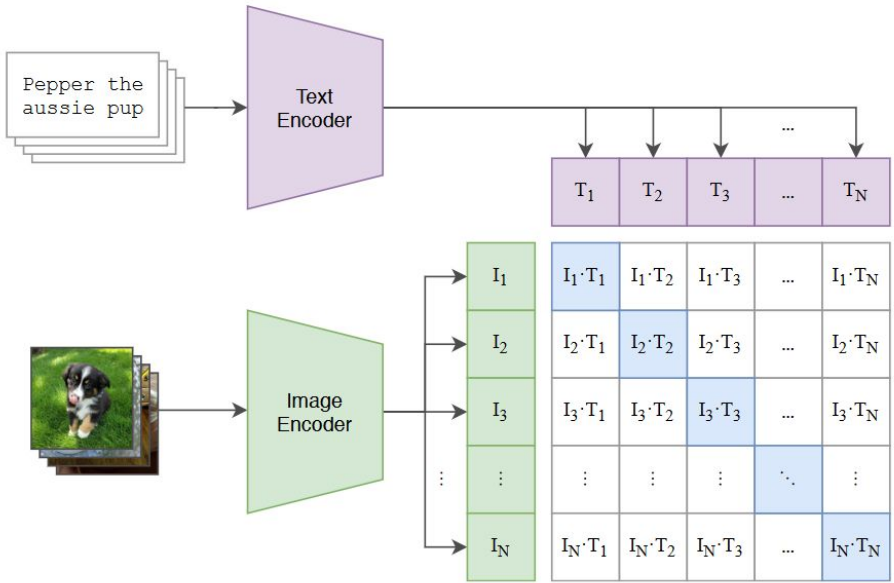
# Key Concepts

•Zero-Shot Transfer: CLIP can perform tasks it was not explicitly trained for without the need for fine-tuning. This is achieved by leveraging the learned association between visual and textual representations

•Contrastive Learning: The model learns by comparing image and text embeddings, maximizing similarity for corresponding pairs and minimizing it for non-corresponding pairs within a batch

•Natural Language Supervision: Text descriptions, rather than manual labels, are used to train the model. This allows the model to learn from a broader range of concepts and contexts

•Image and Text Encoders: CLIP uses a ResNet or Vision Transformer (ViT) for processing images and a Transformer for processing text. Both map their respective inputs into a common embedding space.
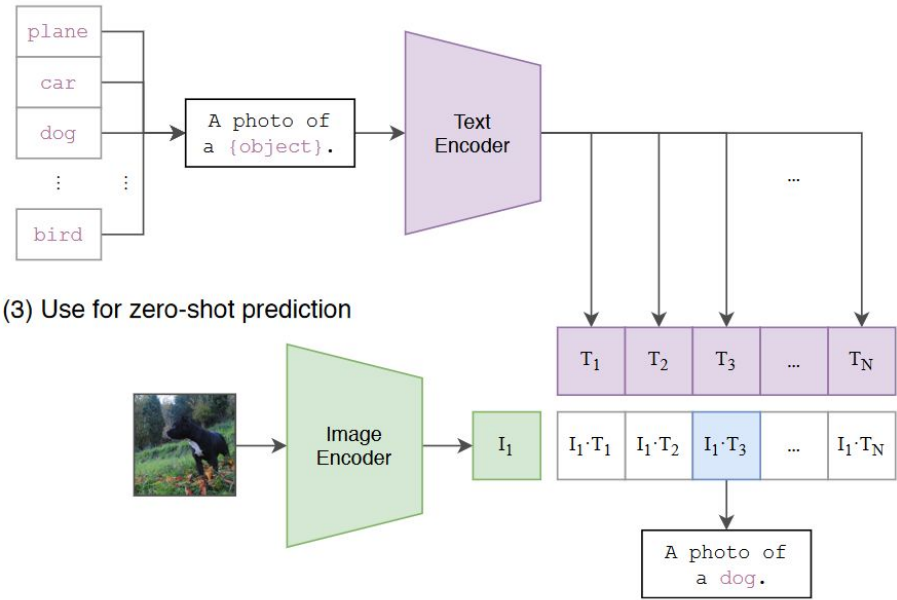
Increase cosine similarity of (image,text) pairs

In inference, it calculates the similarity scores between the vector of a single image with a bunch of possible caption vectors, and picks the caption with the highest similarity.

The paper says they have high zero-shot performance — this is because even though the model might not have been trained on any examples of the classes in the ImageNet dataset, it still performs well because it could kind of figure out what the words of the classes mean and associate that with the images.

# (1) Contrastive pre-training

Pepper the aussie pup → Text Encoder → $T_1$ $T_2$ $T_3$ ... $T_N$

Image Encoder → $I_1$ $I_2$ $I_3$ ... $I_N$

|  | $T_1$ | $T_2$ | $T_3$ | ... | $T_N$ |
|---|---|---|---|---|---|
| $I_1$ | $I_1 \cdot T_1$ | $I_1 \cdot T_2$ | $I_1 \cdot T_3$ | ... | $I_1 \cdot T_N$ |
| $I_2$ | $I_2 \cdot T_1$ | $I_2 \cdot T_2$ | $I_2 \cdot T_3$ | ... | $I_2 \cdot T_N$ |
| $I_3$ | $I_3 \cdot T_1$ | $I_3 \cdot T_2$ | $I_3 \cdot T_3$ | ... | $I_3 \cdot T_N$ |
| ⋮ | ⋮ | ⋮ | ⋮ | ⋱ | ⋮ |
| $I_N$ | $I_N \cdot T_1$ | $I_N \cdot T_2$ | $I_N \cdot T_3$ | ... | $I_N \cdot T_N$ |

# (2) Create dataset classifier from label text

plane
car
dog
⋮
bird

→ A photo of a {object}. → Text Encoder → $T_1$ $T_2$ $T_3$ ... $T_N$

# (3) Use for zero-shot prediction

Image Encoder → $I_1$

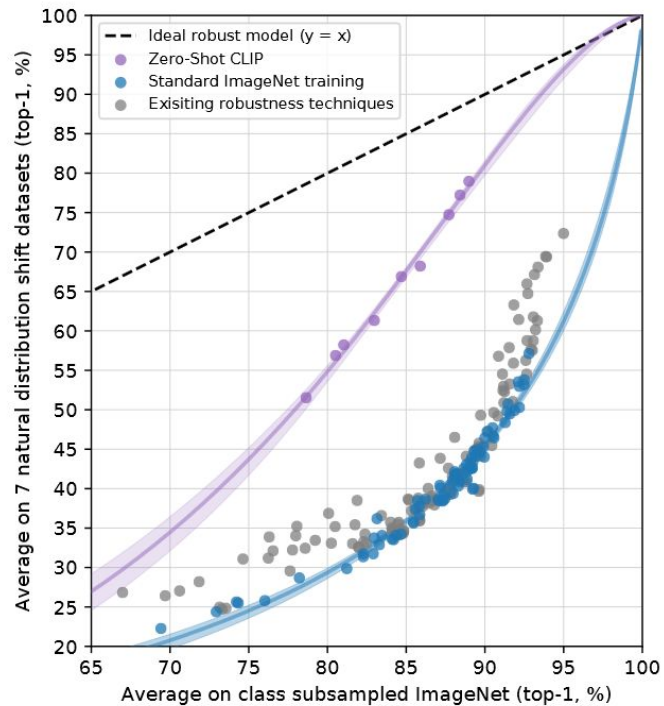|  | $T_1$ | $T_2$ | $T_3$ | ... | $T_N$ |
|---|---|---|---|---|---|
| $I_1$ | $I_1 \cdot T_1$ | $I_1 \cdot T_2$ | $I_1 \cdot T_3$ | ... | $I_1 \cdot T_N$ |

→ A photo of a dog.

# Methodology

• Data Collection: CLIP was trained on a dataset of 400 million (image, text) pairs collected from the internet. This massive dataset facilitates learning a wide variety of visual and textual concepts

• Training Process: The image and text encoders are jointly trained to map images and text to a shared embedding space. A contrastive loss function is used to maximize the cosine similarity between embeddings of corresponding image-text pairs and minimize the similarity of non-corresponding pairs. In each batch, the model attempts to predict the correct text that goes with each image in the batch

• Inference (Zero-Shot): To perform a classification task, the class names are converted into text embeddings, and the input image is converted into an image embedding. The class with the most similar text embedding is then predicted as the output

• Prompt Engineering: The way the text is phrased can significantly affect the model's performance, so using prompts like "a photo of a {label}" provides context and enhances the results

.

# Experimental Results - Zero-Shot Performance

• Broad Applicability: CLIP demonstrates strong performance on a variety of tasks without task-specific training, including image classification, object detection and action recognition

• Comparison to Other Models: CLIP outperforms models trained with traditional methods, such as those trained on ImageNet

• Robustness: CLIP is more robust to distribution shifts compared to models trained on ImageNet, indicating better generalization to new data

• CLIP shows strong zero-shot performance across a wide variety of datasets

| | Dataset Examples | ImageNet ResNet101 | Zero-Shot CLIP | Δ Score |
|---|---|---|---|---|
| ImageNet | | **76.2** | **76.2** | 0% |
| ImageNetV2 | | 64.3 | **70.1** | +5.8% |
| ImageNet-R | | 37.7 | **88.9** | +51.2% |
| ObjectNet | | 32.6 | **72.3** | +39.7% |
| ImageNet Sketch | | 25.2 | **60.2** | +35.0% |
| ImageNet-A | | 2.7 | **77.1** | +74.4% |

# Experimental Results - Linear Probing

• Linear Classification: CLIP's features are effective when used with a linear classifier on downstream tasks, further demonstrating the quality and transferability of the learned representations

• Performance Boost: CLIP shows a significant performance boost when compared to other models in linear probing

• Effectiveness: Linear probing indicates that CLIP's features are highly informative and transferable, showing the model's representation power

# Experimental Results - Bias

•Bias Detection: CLIP reflects biases present in its training data, particularly regarding gender, race, and age

•FairFace Dataset: The authors use the FairFace dataset to explore and quantify the biases present in the model

•Implications: These findings are critical for understanding the limitations and potential ethical concerns for real-world deployment of the model

•Mitigation: The authors demonstrate methods to reduce bias by adjusting thresholds, showing the possibility of mitigating biases

# Comparison to Human Performance

Human-Level Performance: CLIP's performance approaches human-level on some tasks, particularly in tasks that are well-represented in the training data

•Still Below Human Performance: Humans can outperform CLIP in many instances, indicating that there is still a performance gap between the two

•Human Knowledge: Humans can use their existing knowledge to update their prior beliefs, whereas the model does not

•Discrepancies: There is a natural difference in the way humans and CLIP process information that creates differences in performance

|  | Accuracy | Majority Vote on Full Dataset | Accuracy on Guesses | Majority Vote Accuracy on Guesses |
|---|---|---|---|---|
| Zero-shot human | 53.7 | 57.0 | 69.7 | 63.9 |
| Zero-shot CLIP | **93.5** | **93.5** | **93.5** | **93.5** |
| One-shot human | 75.7 | 80.3 | 78.5 | 81.2 |
| Two-shot human | 75.7 | 85.0 | 79.2 | 86.1 |

# Key Findings

•Natural Language Supervision Effectiveness: Large image-text datasets are effective for learning transferable visual representations, which underscores the value of training on such data

•Zero-Shot Transfer: CLIP demonstrates a high degree of adaptability, showing its ability to be applied to a variety of tasks without task-specific training

•Robustness: CLIP demonstrates improved generalization capabilities, being more robust to distribution shifts than ImageNet-trained models

•Bias in Models: CLIP reflects biases inherent in its training data, which is a critical consideration

.

# Limitations

•Zero-Shot Limitations: The paper discusses some limitations of zero-shot transfer, and situations where fine tuning the model would be more effective

•Potential Issues: There is a continued need for model improvement, especially in situations where the model's performance is below optimal

•Ongoing Research: There are opportunities for improving the model's performance on several tasks, which suggests an area for future research

# Additional Notes

•Ablation Study: The paper includes an ablation study of different CLIP components, providing insights into their individual impact on performance

•Model Details: The paper provides detailed information about the model's architecture, training process, and datasets used

•Evaluation: The model is evaluated on 27 datasets, demonstrating its ability to be generalized to a wide range of tasks

•Other Tasks: CLIP is applicable to other tasks beyond image classification, such as image-text retrieval and OCR

•Societal Impact: The paper discusses the biases in the model and their societal impact, which is crucial for the ethical consideration of AI models

# Conclusion

- CLIP is a powerful model that leverages natural language supervision to learn transferable visual representations

- The model shows impressive zero-shot transfer capabilities and robustness to distribution shift

- It is important to consider the biases inherent in the model

- Future research can explore the societal impacts of AI models

# Future Work

•Improve Model Robustness: Further efforts are needed to make the model more robust to distribution shifts

•Mitigate Bias: It's crucial to explore methods for minimizing biases in the model

•Societal Impacts: Further exploration of the societal impacts of AI models is essential

•Future research should also focus on improving zero-shot transfer capabilities

# Relation to Captcha Task

The task should be very easy to solve while using Clip. I tried implementing it but face limitations in execute time with gpu in colab and limited resources in my PC.