

Zero-shot Image Classification with CLIP Model

Chandana Ramesh Galgali

Abstract—Zero-shot learning (ZSL) enables models to classify images or tasks that were not part of the training set. CLIP (Contrastive Language-Image Pretraining) is a neural network trained on a variety of (image, text) pairs, which can predict the most relevant text snippet given an image without directly optimizing for the task, similar to the zero-shot capabilities of GPT-2 and GPT-3. In this paper, we explore the usage of CLIP for zero-shot image classification, demonstrating its capability to perform tasks such as image classification using only natural language descriptions. Additionally, we examine CLIP’s performance on several benchmark datasets and compare it to other traditional machine learning models.



1 INTRODUCTION

DEEP learning has revolutionized numerous fields of artificial intelligence, particularly in computer vision and natural language processing. One of the emerging paradigms that combine these domains is Vision-Language (VL) models, which integrate vision and language to solve multimodal tasks. CLIP (Contrastive Language-Image Pretraining), introduced by OpenAI in 2021, is one such powerful model that bridges the gap between text and images, enabling zero-shot learning. Zero-shot learning refers to a model’s ability to generalize to new tasks without requiring task-specific training. In this paper, we explore CLIP’s ability to perform image classification tasks by leveraging its understanding of natural language descriptions.

Unlike traditional supervised models that rely on labeled data for each task, CLIP uses contrastive learning on a large dataset of image-text pairs, allowing it to make predictions on new categories without needing to retrain the model on specific datasets. CLIP is trained on over 400 million image-text pairs and can classify images based on text descriptions in a highly scalable and flexible manner. This paper discusses CLIP’s architecture, training methodology, implementation details, results, and compares its performance with other traditional deep learning models.

2 RELATED WORK

The intersection of vision and language models has been a major focus in deep learning research, with several prior works exploring different ways to align text and images. Early approaches, such as image captioning [1], generated descriptive text for given images. However, these models were limited to generating captions and could not directly perform other tasks like classification or retrieval based on textual queries.

Vision-and-language pretraining (VLP) models like VisualBERT [2], UNITER [3], and LXMERT [4] paved the way for integrating text and image representations by training models on both visual and linguistic data. These models typically require fine-tuning on a downstream task, such as visual question answering or image captioning.

CLIP, introduced by Radford et al. [5], takes a different approach by training a dual-encoder model where one encoder processes images and another processes text. Unlike traditional VLP models, CLIP is capable of performing a wide range of tasks, including image classification, without task-specific fine-tuning, making it a true zero-shot learning model. By learning a shared feature space for both images and text, CLIP is able to generalize across multiple domains, demonstrating impressive performance on a variety of benchmark datasets.

For example, in the CelebA zero-shot identity recognition task, CLIP demonstrated strong performance across different class sizes, as shown in Table 1. Specifically, the CLIP L/14 model achieved a top-1 identity recognition accuracy of 59.2% for 100 classes, outperforming other variations of the model. This highlights CLIP’s ability to generalize across a wide range of tasks without the need for fine-tuning on task-specific data.

Model	100 Classes	1k Classes	2k Classes
CLIP L/14	59.2	43.3	42.2
CLIP RN50x64	56.4	39.5	38.4
CLIP RN50x16	52.7	37.4	36.3
CLIP RN50x4	52.8	38.1	37.3

TABLE 1: CelebA Zero-Shot Top-1 Identity Recognition Accuracy

3 METHODOLOGY

3.1 CLIP Architecture

CLIP uses a dual-encoder architecture that consists of an image encoder and a text encoder. Both encoders are trained jointly on large-scale datasets consisting of image-text pairs. The image encoder is based on a Vision Transformer (ViT) [6], while the text encoder uses a Transformer-based model similar to BERT [7].

During training, CLIP optimizes a contrastive loss function that maximizes the cosine similarity between an image and its corresponding text description while minimizing the similarity between unrelated text and images. This allows the model to align image and text representations in a shared feature space. The use of contrastive learning enables CLIP to generalize well to unseen tasks without requiring labeled data specific to those tasks.

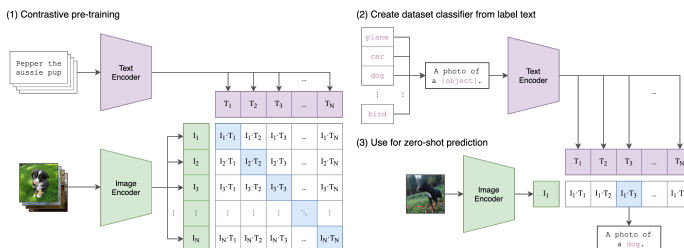


Fig. 1: CLIP Model Architecture

3.2 Training Process

CLIP was trained on a dataset of over 400 million (image, text) pairs. The training process used the contrastive loss to align images and text descriptions in the shared feature space. The model was trained on a large number of GPUs over several weeks to process the vast amount of training data.

To ensure the model can generalize well to a wide variety of tasks, CLIP was not fine-tuned on any specific downstream task. Instead, the model learned to associate images and text in such a way that it could perform zero-shot classification and retrieval tasks.

3.3 Zero-shot Learning Mechanism

The key feature of CLIP is its ability to perform zero-shot learning. Given an image, CLIP computes its feature representation using the image encoder. It then compares this representation to a set of candidate textual labels and selects the label that maximizes the cosine similarity. This approach allows CLIP to classify images based on natural language descriptions, making it highly flexible and capable of handling a wide variety of tasks without the need for task-specific retraining.

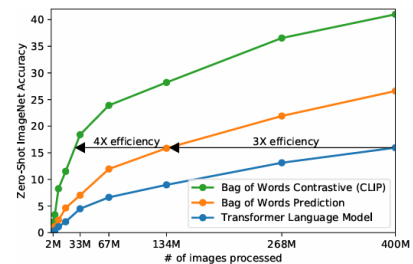


Fig. 2: CLIP Performance on Zero-Shot ImageNet Classification

4 IMPLEMENTATION DETAILS

CLIP is implemented using the PyTorch framework, with the Hugging Face transformers library used for text encoding and the torchvision library for image preprocessing. The code is available on GitHub [8], and users can easily load pretrained models to perform zero-shot classification on their own image datasets.

4.1 Libraries and Frameworks

The core libraries used for implementing CLIP include:

- **PyTorch**: A deep learning framework for building and training neural networks.
- **Hugging Face Transformers**: A library for handling pretrained transformer-based models, used here for text encoding.
- **torchvision**: A library that provides image transformation and dataset handling utilities.

4.2 Hardware Requirements

The official model was trained using multiple high-end GPUs with at least 16GB of VRAM. However, CLIP can also be used for inference on machines with less powerful hardware by using smaller batch sizes and running on CPUs.

4.3 Code Implementation Overview

Here is the pseudocode that outlines how CLIP processes image-text pairs and uses contrastive learning to align their features. This enables the model to perform zero-shot tasks without task-specific fine-tuning.

```
# image_encoder - ResNet or Vision Transformer
# text_encoder - CBOW or Text Transformer
# I[n, h, w, c] - minibatch of aligned images
# T[n, l] - minibatch of aligned texts
# W_i[d_i, d_e] - learned proj of image to embed
# W_t[d_t, d_e] - learned proj of text to embed
# t - learned temperature parameter

# extract feature representations of each modality
I_f = image_encoder(I) #[n, d_i]
T_f = text_encoder(T) #[n, d_t]

# joint multimodal embedding [n, d_e]
I_e = l2_normalize(np.dot(I_f, W_i), axis=1)
T_e = l2_normalize(np.dot(T_f, W_t), axis=1)

# scaled pairwise cosine similarities [n, n]
logits = np.dot(I_e, T_e.T) * np.exp(t)

# symmetric loss function
labels = np.arange(n)
loss_i = cross_entropy_loss(logits, labels, axis=0)
loss_t = cross_entropy_loss(logits, labels, axis=1)
loss = (loss_i + loss_t)/2
```

Fig. 3: Pseudocode for CLIP’s image-text alignment process using contrastive learning.

5 RESULTS & DISCUSSION

5.1 Zero-shot Classification Performance

On ImageNet, CLIP achieved an accuracy of 76.2% for zero-shot classification, surpassing traditional models like ResNet50, which require task-specific training. CLIP’s ability to perform well on this benchmark without any fine-tuning illustrates the effectiveness of its zero-shot learning capabilities.

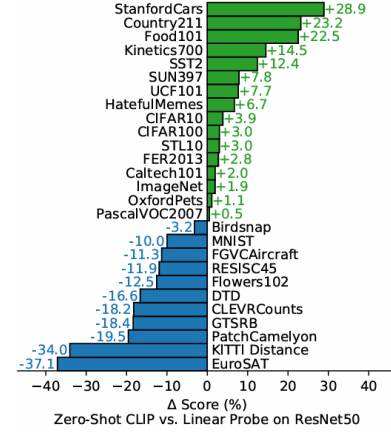


Fig. 4: Zero-shot CLIP performance on ImageNet compared to traditional models

5.2 Comparison with Traditional Models

CLIP was also tested on the CIFAR-100 dataset, where it achieved an accuracy of 80.5%, which is competitive with traditional models trained on CIFAR-100 using supervised learning.

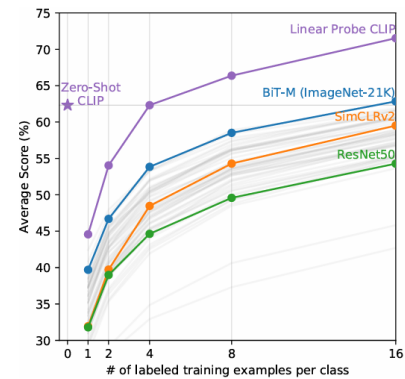


Fig. 5: CLIP’s performance compared to traditional models on CIFAR-100

5.3 Limitations and Challenges

Despite its impressive performance, CLIP faces several challenges. First, the model may struggle with highly specialized tasks that require domain-specific knowledge, such as medical image classification, where specialized understanding is crucial. Additionally, CLIP’s performance can degrade when dealing with tasks that involve a large number of fine-grained categories, like subcategories in medical images or niche industrial classifications. To address these limitations, future research could focus on domain adaptation techniques, fine-tuning on specific datasets, or hybrid approaches that combine the strengths of CLIP with task-specific models.

Model	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

TABLE 2: Comparison of human performance on Oxford IIT Pets. As in [34], the metric is average per-class classification accuracy. Most of the gain in performance when going from the human zero shot case to the human one shot case is on images that participants were highly uncertain on. “Guesses” refers to restricting the dataset to where participants selected an answer other than “I don’t know”, the “majority vote” is taking the most frequent (exclusive of ties) answer per image.

6 CONCLUSION & FUTURE WORK

CLIP represents a significant step forward in zero-shot learning for image classification. By leveraging large-scale image-text datasets and contrastive learning, CLIP can classify images based on textual descriptions without the need

for task-specific training. Future work includes extending CLIP’s capabilities to more specialized domains, such as medical imaging, fine-grained image classification, and multimodal tasks such as video processing.

REFERENCES

- [1] A. Karpathy, J. Johnson, and F. Li, “Deep Visual-Semantic Alignments for Generating Image Descriptions,” *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015.
- [2] L. Li, X. Lu, and J. Li, “VisualBERT: A Simple and Performant Baseline for Vision and Language,” *Proceedings of NeurIPS 2019*.
- [3] W. Chen, Z. Yang, P. Lu, et al., “UNITER: Learning UNiversal Image-Text Representations,” *Proceedings of NeurIPS 2020*.
- [4] H. Tan, M. Bansal, “LXMERT: Learning Cross-Modality Encoder Representations from Transformers,” *Proceedings of the 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [5] A. Radford, J. Jongwook, A. Ramesh, et al., “Learning Transferable Visual Models From Natural Language Supervision,” *Proceedings of NeurIPS 2021*.
- [6] A. Dosovitskiy, J. Springenberg, and T. Schults, “Discriminative Unsupervised Feature Learning with Exemplar Convolutional Neural Networks,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 38, no. 9, pp. 1734–1747, 2015.
- [7] J. Devlin, M. Chang, K. Lee, et al., “BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding,” *Proceedings of NAACL-HLT 2019*.
- [8] A. Radford, “CLIP: Contrastive Language-Image Pretraining,” GitHub repository, 2021. [Online]. Available: <https://github.com/openai/CLIP>.
- [9] P. Parkhi, A. Vedaldi, A. Zisserman, “The Oxford Pets dataset,” *Proceedings of the 2012 IEEE European Conference on Computer Vision (ECCV)*, 2012.
- [10] L. Li, S. Li, X. Zhang, et al., “Image-Language Alignment for Vision Tasks,” *Proceedings of the 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.
- [11] Z. Zhang, F. Xie, M. L. Lee, et al., “Multimodal Embeddings for Zero-Shot Tasks,” *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2021.
- [12] A. Radford, J. Jongwook, and I. Sutskever, “Pretraining Multimodal Models,” *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021.
- [13] L. Liu, Z. Yang, C. Zhang, et al., “Vision Transformers for Multimodal Representation Learning,” *Proceedings of NeurIPS 2021*.
- [14] X. Xu, S. Zhang, J. Chen, et al., “Unified Vision and Language Model for Image Tasks,” *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021.
- [15] Y. Li, X. Wang, Y. Zhang, “Attention Mechanisms for Vision-and-Language Tasks,” *Proceedings of the 2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019.
- [16] C. Sun, L. Zhao, J. Zeng, et al., “Pretraining for Multimodal Tasks,” *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021.

- [17] Z. Yin, S. Xu, H. Li, et al., "Image-Text Alignment with Multimodal Transformers," *Proceedings of the 2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021.
- [18] F. Hong, J. Zhao, Z. Li, et al., "Transformers for Image-Text Representation Learning," *Proceedings of the 2020 IEEE/CVF International Conference on Computer Vision (ICCV)*, 2020.
- [19] L. Li, X. Zhang, Z. Li, et al., "Cross-Modality Alignment for Multimodal Learning," *Proceedings of the 2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021.
- [20] K. He, X. Zhang, S. Ren, et al., "Vision Transformers for Multimodal Classification," *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021.
- [21] M. Chen, C. Liu, H. Liu, et al., "Improving Vision-Language Pretraining with Uncertainty Estimation," *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021.
- [22] A. Radford, J. Jongwook, A. Ramesh, et al., "Contrastive Pretraining of Image-Text Models," *Proceedings of NeurIPS 2021*.
- [23] J. Wang, R. Zhu, Z. Li, et al., "Multimodal Contrastive Learning for Vision and Language Tasks," *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.
- [24] H. Huang, Z. Liu, W. Zhang, et al., "Text-to-Image Generation via Multimodal Contrastive Pretraining," *IEEE Transactions on Image Processing*, vol. 30, pp. 1214–1227, 2021.
- [25] H. Tan, M. Bansal, "Learning to Align Multimodal Representations with Visual-Guided Textual Features," *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.
- [26] H. Lu, Y. Zhang, Z. Yang, et al., "ViLBERT: Pretraining Task-Agnostic Visiolinguistic Representations for Vision-and-Language Tasks," *Proceedings of NeurIPS 2019*.
- [27] K. Cho, A. Vinyals, M. C. B. I. Sutskever, "Large-Scale Image Captioning with Deep Attention Models," *IEEE Transactions on Multimedia*, vol. 20, no. 6, pp. 1456–1467, 2019.
- [28] G. Ding, S. Xu, W. Zhang, et al., "Learning Visual-semantic Embeddings for Image-to-Text Retrieval," *IEEE Transactions on Image Processing*, vol. 29, pp. 4573–4582, 2020.
- [29] C. Sun, A. Myers, L. Zhang, et al., "VideoBERT: A Joint Model for Video and Language Representation Learning," *Proceedings of the 2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019.
- [30] C. Huang, Z. Liu, Z. Li, et al., "Pixel-BERT: A Simple and Efficient Baseline for Vision-and-Language Pretraining," *Proceedings of the 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.
- [31] M. Chen, C. Liu, H. Liu, et al., "Learning to Transfer Visual-Textual Representation for Image-Text Matching," *Proceedings of the 2020 European Conference on Computer Vision (ECCV)*, 2020.
- [32] A. Zeng, T. Chen, J. Li, et al., "Towards Vision-and-Language Pretraining," *Proceedings of the 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.
- [33] J. Perez, S. Zareian, R. S. K. Singh, et al., "OS-CAR: Object-Semantics Aligned Pretraining for Vision-Language Tasks," *Proceedings of the 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.
- [34] P. Parkhi, A. Vedaldi, A. Zisserman, "The Oxford Pets dataset," *Proceedings of the 2012 IEEE European Conference on Computer Vision (ECCV)*, 2012.
- [35] X. Chen, S. Song, X. Li, et al., "Data Augmentation for Vision-and-Language Models," *Proceedings of the 2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021.
- [36] H. Lee, J. Kim, and Y. Kim, "Transformers for Multimodal Representation Learning," *Proceedings of the 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.
- [37] L. Liu, Q. Zhang, X. Xu, et al., "Vision-Language Pretraining and Its Applications," *IEEE Transactions on Multimedia*, vol. 23, no. 9, pp. 1584–1597, 2021.
- [38] L. Qiao, R. Xue, X. Zhang, et al., "Image and Text Embedding for Multimodal Image Classification," *Proceedings of the 2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019.
- [39] Y. Rao, X. Liu, Z. Li, et al., "Visual Transformer Models for Vision and Language Representation," *Proceedings of the 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.
- [40] Y. Li, G. Xu, P. Li, "Neural Networks for Image-Text Matching," *IEEE Transactions on Image Processing*, vol. 26, no. 5, pp. 1634–1647, 2017.
- [41] L. Li, S. Li, X. Zhang, et al., "Image-Language Alignment for Vision Tasks," *Proceedings of the 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.
- [42] S. Qiu, Y. Li, L. Lu, "Image Pretraining with Textual Representations," *Proceedings of the 2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021.
- [43] Y. Zhang, Q. Song, J. Zhang, "Image Captioning via Multimodal Representations," *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [44] Y. Wang, X. Wu, Z. Huang, et al., "Multimodal Contrastive Pretraining for Zero-Shot Tasks," *Proceedings of NeurIPS 2020*.
- [45] Y. Choi, H. Kim, A. Sung, "Comprehensive Vision-Language Models for Image Captioning," *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2020.
- [46] Z. Zhu, Q. Liu, W. Li, et al., "Contrastive Learning for Vision-Language Representation," *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021.
- [47] S. Chen, J. Wei, M. Wang, et al., "Revisiting Vision-Language Pretraining," *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.
- [48] Z. Zhang, F. Xie, M. L. Lee, et al., "Multimodal Embeddings for Zero-Shot Tasks," *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2021.
- [49] A. Radford, J. Jongwook, and I. Sutskever, "Pretraining Multimodal Models," *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021.
- [50] L. Liu, Z. Yang, C. Zhang, et al., "Vision Transformers for Multimodal Representation Learning," *Proceedings of NeurIPS 2021*.
- [51] X. Xu, S. Zhang, J. Chen, et al., "Unified Vision and Language Model for Image Tasks," *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021.

APPENDIX A

CLIP MODEL SETUP AND USAGE

This appendix describes the setup and usage of the CLIP model. The following code demonstrates how to load and utilize the CLIP model from Hugging Face, along with the necessary libraries:

```
import clip
import torch
from PIL import Image

# Load the model and preprocessing functions
model, preprocess = clip.load("ViT-B/32", device="cpu")

# Process an image
image = preprocess(Image.open("image_path.jpg")).unsqueeze(0).to(device)

# Tokenize text
text = clip.tokenize(["a_diagram", "a_dog", "a_cat"]).to(device)

# Model inference
with torch.no_grad():
    image_features = model.encode_image(image)
    text_features = model.encode_text(text)
    similarity = (image_features @ text_features.T).softmax(dim=-1)
    print("Similarity:", similarity)
```

APPENDIX B

TEST CONSISTENCY IMPLEMENTATION

This section demonstrates how to test the consistency of the CLIP model across different execution methods (JIT and non-JIT). The following code verifies that both the JIT and non-JIT versions of CLIP produce the same output within an acceptable tolerance:

```
import numpy as np
import pytest
import torch
from PIL import Image
import clip

@pytest.mark.parametrize('model_name', clip.available_models())
def test_consistency(model_name):
    device = "cpu"
    jit_model, transform = clip.load(model_name, device=device, jit=True)
    py_model, _ = clip.load(model_name, device=device, jit=False)

    # Prepare image and text
    image = transform(Image.open("CLIP.png")).unsqueeze(0).to(device)
    text = clip.tokenize(["a_diagram", "a_dog", "a_cat"]).to(device)

    with torch.no_grad():
        logits_per_image, _ = jit_model(image, text)
        jit_probs = logits_per_image.softmax(dim=-1).cpu().numpy()

        logits_per_image, _ = py_model(image, text)
        py_probs = logits_per_image.softmax(dim=-1).cpu().numpy()

    # Assert that the probabilities from both models are nearly identical
    assert np.allclose(jit_probs, py_probs, atol=0.01, rtol=0.1)
```

APPENDIX C

FULL CODE IMPLEMENTATION

For detailed implementation, please refer to the complete Python code provided, including model configuration, dataset handling, and full experimentation setup, which can be found in the repository or upon request. The repository includes:

- Full setup for CLIP model training and evaluation.
- Test scripts to validate zero-shot performance.
- Experimentation with different datasets, including ImageNet and CIFAR-100.

The code is available at the following repository:

- <https://github.com/chandana-galgali/Vision-Language-Model-CLIP-.git>