

CLIP

Learning Transferable Visual Models From Natural Language Supervision (2021) 6.1k

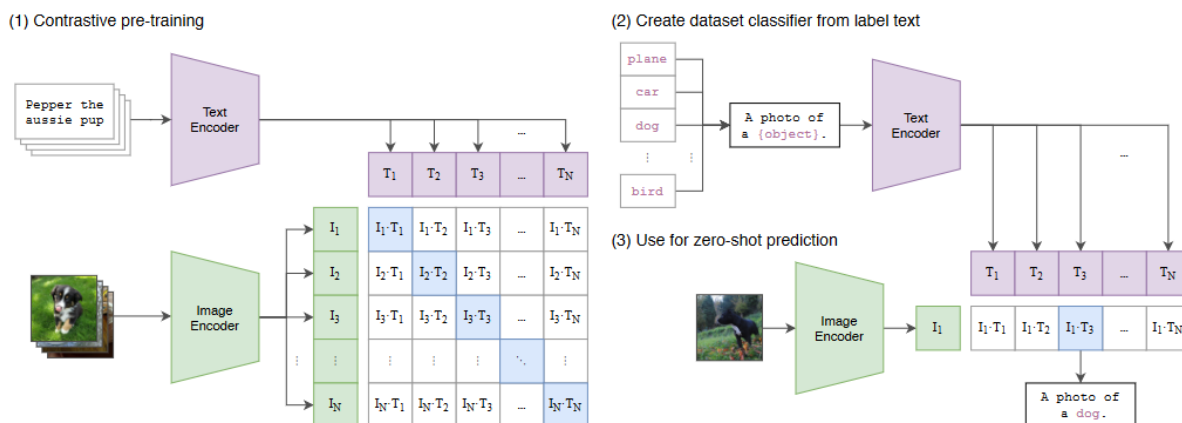


Figure 1. Summary of our approach. While standard image models jointly train an image feature extractor and a linear classifier to predict some label, CLIP jointly trains an image encoder and a text encoder to predict the correct pairings of a batch of (image, text) training examples. At test time the learned text encoder synthesizes a zero-shot linear classifier by embedding the names or descriptions of the target dataset's classes.

fixed set of predetermined object categories → limited generality

CLIP: Contrastive Language-Image Pre-training, no classifier or categorical label, transferable to new labels / datasets / tasks: e.g., zero-shot accuracy on ImageNet = ResNet-50, more robust

Contrastive pre-training

larger dataset of (image, text) pairs, larger model

due to training efficiency, use contrastive methods instead of gpt; prediction task is too hard, contrastive task is more reasonable

image_encoder(I), text_encoder(T) → cosine similarities, positive samples on diagonal entry

Zero-Shot Transfer

pretrained text_encoder: labels → text features for all label →

pretrained image_encoder: image → image feature → cosine similarity → softmax

Prompt Engineering and Ensembling

PS. a good research topic in fine-tuning / reasoning stage that doesn't require too much computational resources

Issues:

- polysemy: multiple meanings of the same word

- in pre-training dataset, text = full sentence → in classification task, single word
 - A photo of a {label}. accuracy +1.3%
 - A bright photo of a {label}, a type of pet. 80 templates

Representation Learning (using all data in downstream tasks)

Common methods:

- linear probe: fix the model, fit a linear classifier on a representation extracted from the model and measure its performance
- fine tune: measure the performance of end-to-end fine-tuning of the model

Limitation

Poor on

- fine-grained classification such as differentiating models of cars, species of flowers
- more abstract and systematic tasks such as counting the number of objects in an image.
- novel tasks which are unlikely to be included in CLIP's pre-training dataset, such as classifying

the distance to the nearest car in a photo

- data that is truly out-of-distribution for it, only 88% accuracy on MNIST

Repeatedly query performance on full validation sets to guide the development of CLIP → not truly zero-shot transfer

Related Papers

Segmentation

- Lseg: **Language-driven Semantic Segmentation** (2022) 135
- GroupViT: **GroupViT: Semantic Segmentation Emerges from Text Supervision** (2022) 121

Detection

- ViLD: **Open-vocabulary Object Detection via Vision and Language Knowledge Distillation** (2021) 299
- GLIP: **Grounded Language-Image Pre-training** (2021) 225

Graphics

- CLIPasso: **Semantically-Aware Object Sketching** (2022) 48

Video

- CLIP4Clip: **An Empirical Study of CLIP for End to End Video Clip Retrieval** (2021) 295
- ActionCLIP: **A New Paradigm for Video Action Recognition** (2021) 118

CLIPasso: Semantically-Aware Object Sketching

An object sketching method that can achieve different levels of abstraction

Bezier curve: use four points to control a curve

Loss between the embeddings of the sketch $CLIP(R(\{s_i\}_{i=1}^n))$ and image $CLIP(I)$:

$L_{semantic} = dist(CLIP(I), CLIP(R(\{s_i\}_{i=1}^n)))$ for high-level semantic attributes

$L_{semantic} = \|CLIP_l(I) - CLIP_l(R(\{s_i\}_{i=1}^n))\|_2^2$ on intermediate level activations of CLIP, for low-level spatial features