# OpenAI

# CLIP
Learning Transferable Visual Models From Natural Language Supervision

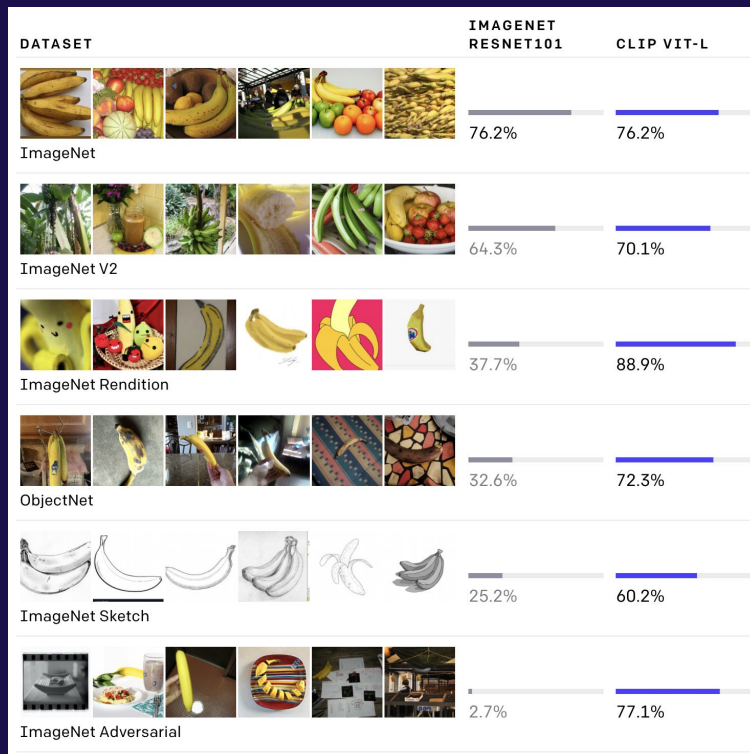Alec Radford, Jong Wook Kim, et al.
January 2021

# What makes CLIP special?

Motivation:

Instead of using a fixed set of labels,
Get supervision from natural language

Result:

Robust zero-shot inference
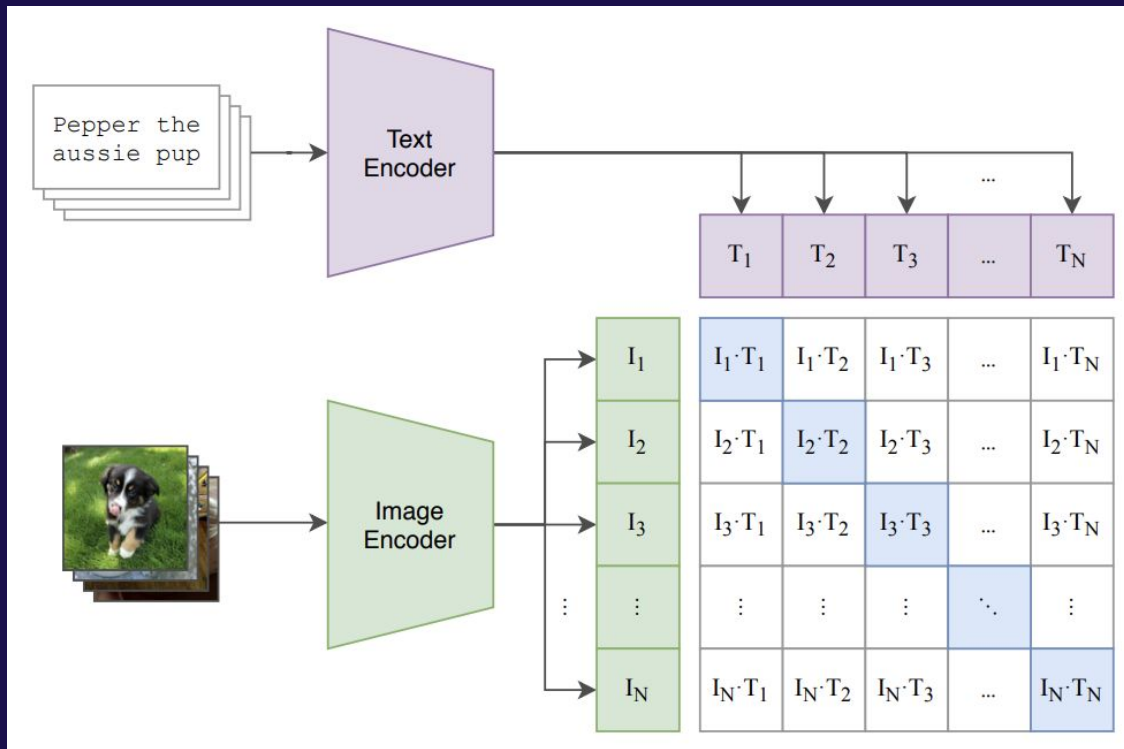Multimodal embedding space

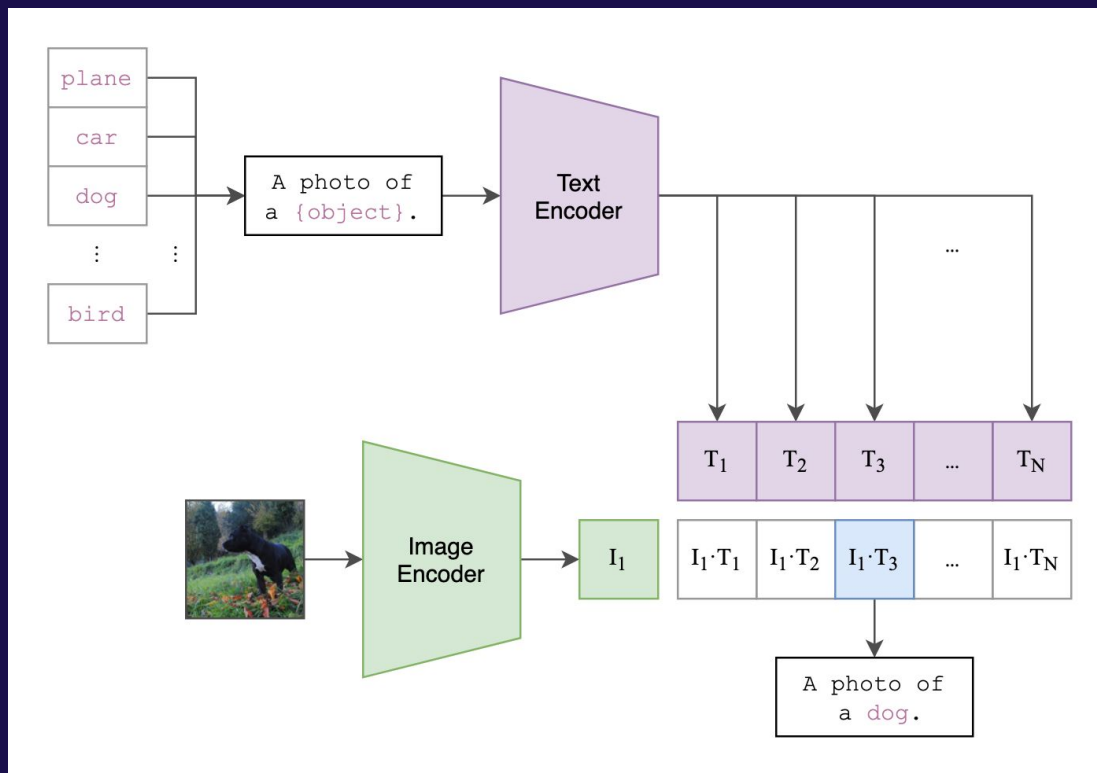| DATASET | IMAGENET RESNET101 | CLIP VIT-L |
|---|---|---|
| ImageNet | 76.2% | 76.2% |
| ImageNet V2 | 64.3% | 70.1% |
| ImageNet Rendition | 37.7% | 88.9% |
| ObjectNet | 32.6% | 72.3% |
| ImageNet Sketch | 25.2% | 60.2% |
| ImageNet Adversarial | 2.7% | 77.1% |

# How does it work?



Pig        Tiger        Panda        Hippo        Camel

# CLIP: Contrastive Language-Image Pre-training

# Zero-shot image classification

# Why contrastive



Zero-shot ImageNet accuracy

40%

30%

20%

10%

0%

Bag of Words Contrastive (CLIP)

Bag of Words Prediction

Transformer Language Model

4x efficiency

3x efficiency

2M  33M  67M  134M  268M  400M

Images processed

**Some CLIP details**

Training
- Trained on 400M image-text pairs from the internet
- Batch size of 32,768
- 32 epochs over the dataset
- Cosine learning rate decay

Architecture
- ResNet-based or ViT-based image encoder
- Transformer-based text encoder

# Representation Learning

## Linear probe
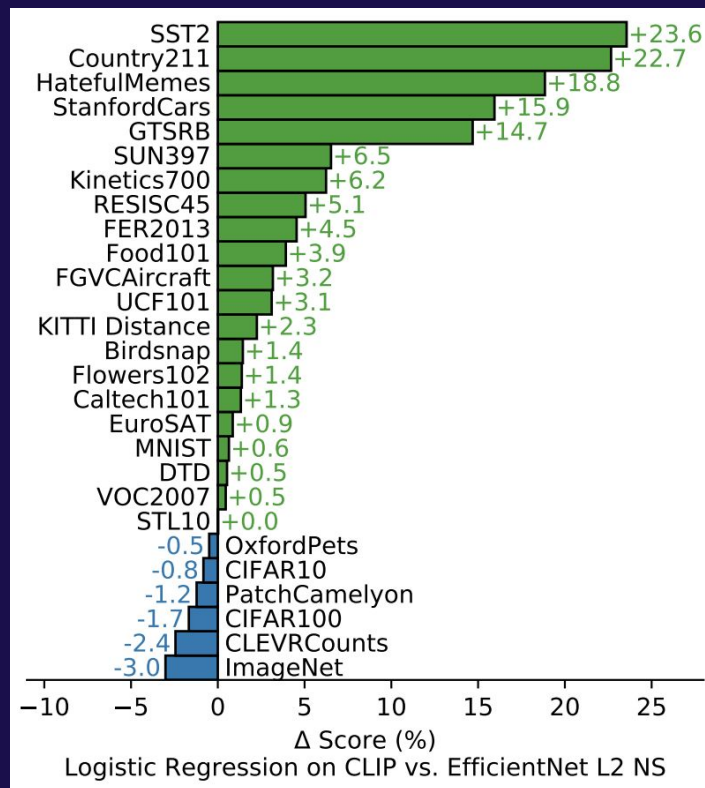
Logistic regression classifier on image features

- L-BFGS
- Only one hyperparameter
- Allows "fair" comparisons with other vision models
- Provides lower bound for fine-tuned models

Evaluated on 27 datasets × 65 vision models

# Linear probe performance vs SOTA vision models



Linear probe average over all 27 datasets

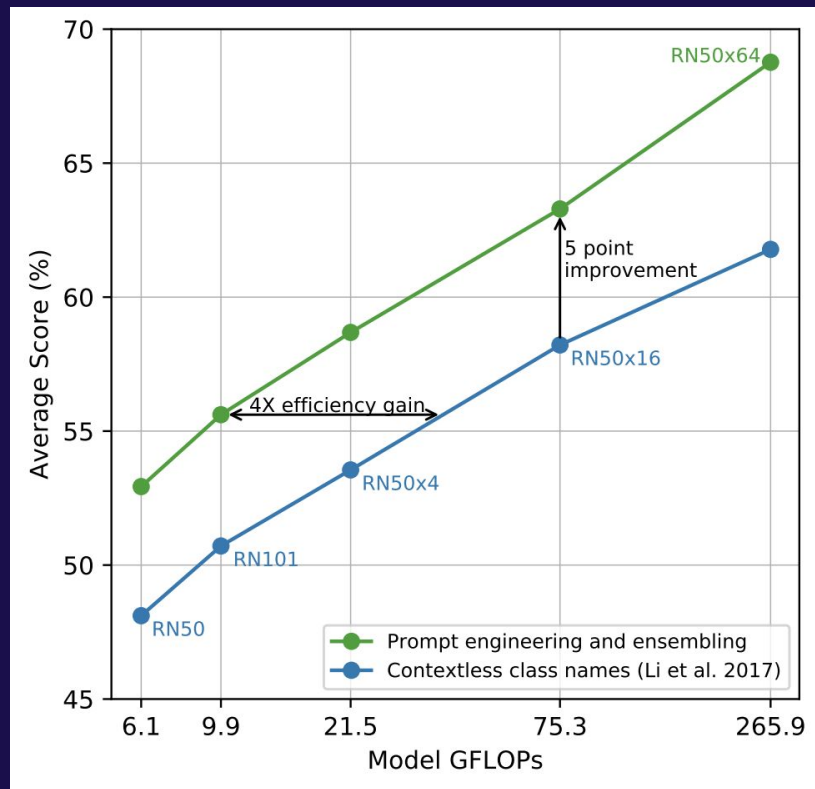# Linear-probe CLIP vs Linear-probe EfficientNet-L2

# vs ImageNet score



Linear probe average over 26 datasets

# Zero-Shot Transfer

# Prompt engineering

# Zero-shot vs Linear-probe ResNet-50
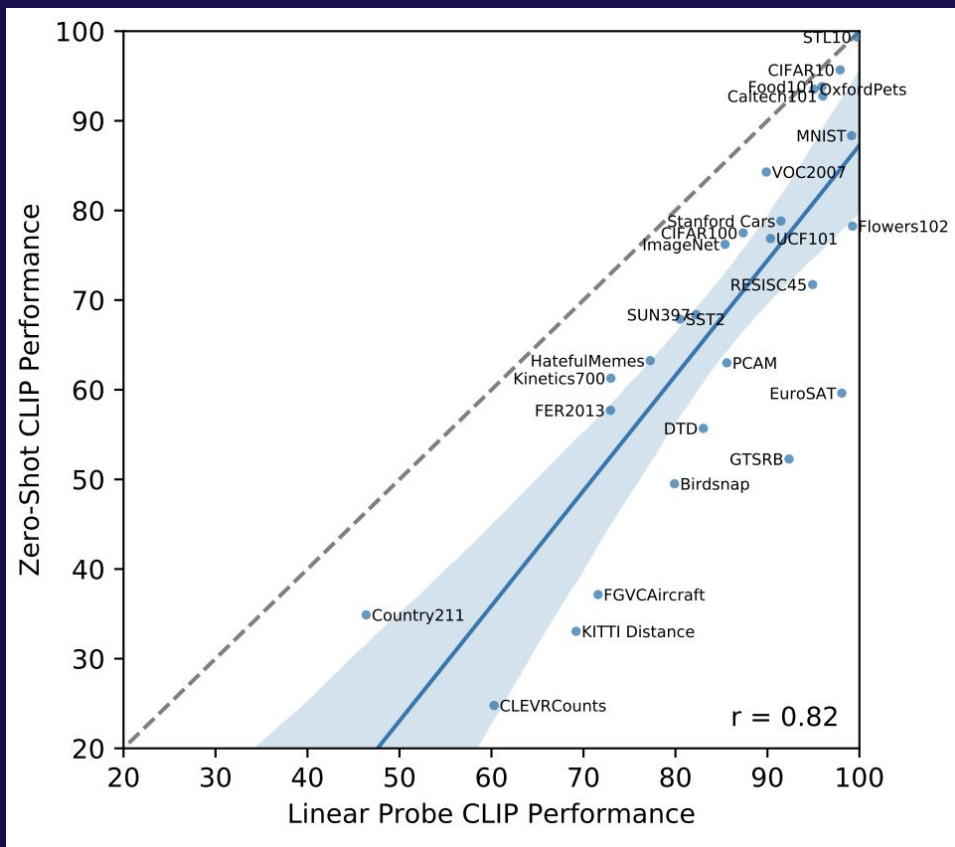
Zero-shot CLIP outperforms ResNet-50 on 16 of 27 datasets



Δ Score (%)
Zero-Shot CLIP vs. Linear Probe on ResNet50

# Zero-shot CLIP vs Few-shot linear probes

Zero-shot CLIP is as good as

- 4-shot linear-probe CLIP
- 16-shot BiT-M

# Zero-shot vs Linear-probe CLIP

# Zero-shot performance vs model size

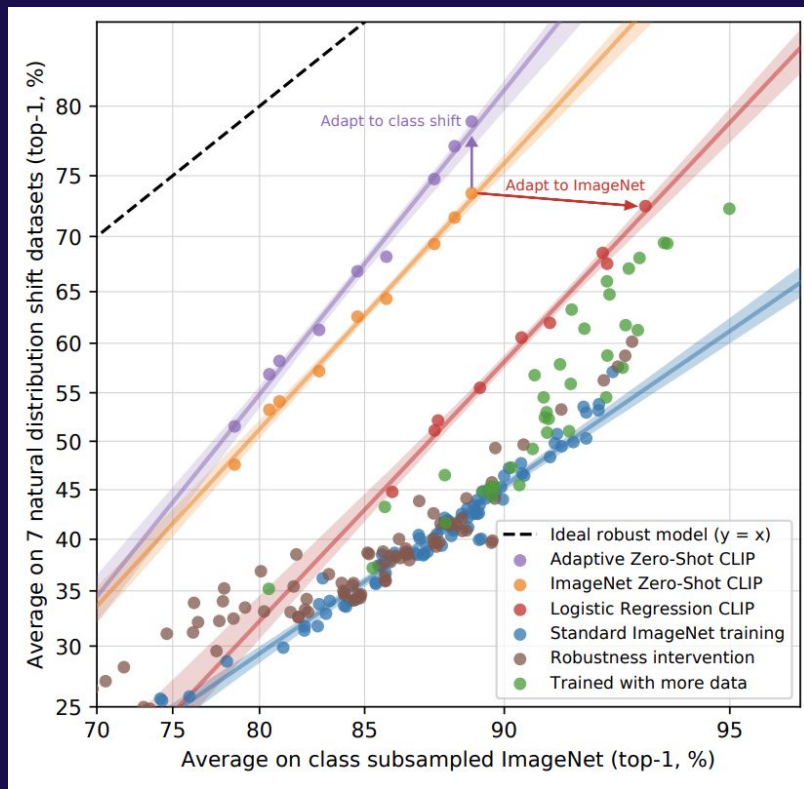# Robustness to Natural Distribution Shift

# Robustness to natural distribution shift

CLIP is significantly more robust!
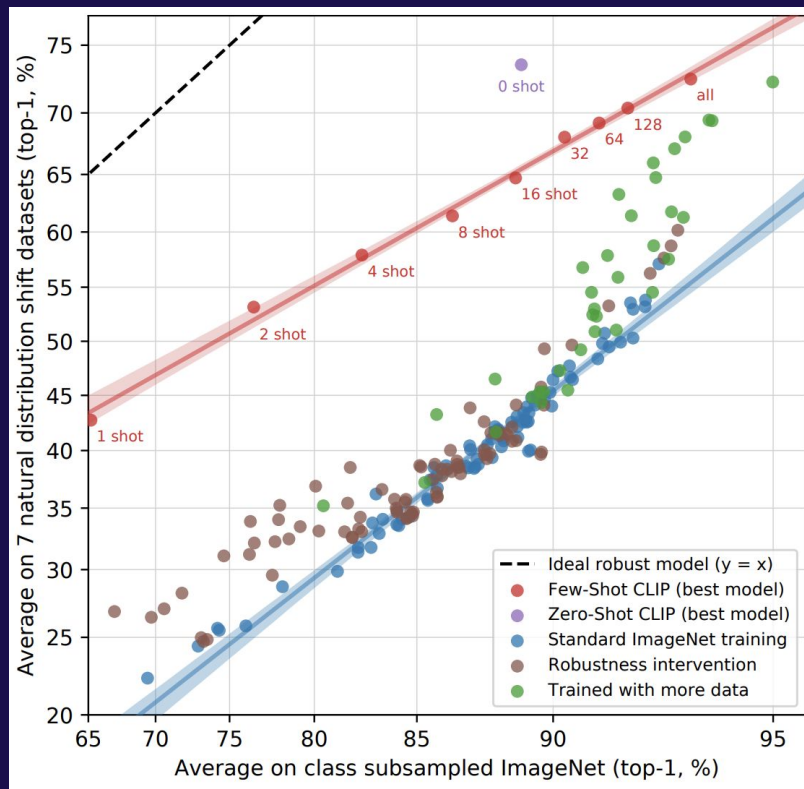
7 ImageNet-like Datasets (Taori et al.)
- ImageNetV2
- ImageNet-A
- ImageNet-R
- ImageNet Sketch
- ObjectNet
- ImageNet Vid
- Youtube-BB

# Adapting to ImageNet does not help robustness

# Robustness of few-shot linear probes

# Limitations and Broader Impacts

**Limitations of CLIP**

- Zero-shot performance is well below the SOTA

- Especially weak on abstract tasks such as counting

- Poor on out-of-distribution data such as MNIST

- Susceptible to adversarial attacks

- Dataset selection bias

- Social biases

**Quantifying the (un)safety of CLIP models**

CLIP has societal biases
- Race
- Gender
- Age

Surveillance usage
- Zero-shot scene classification
- Zero-shot identification of celebrities

Not comprehensive, will continue research to ensure safety
Model card limits usage of CLIP to research-only

# Related Work

Multimodal learning
- VirTex
- ICMLM
- ConVIRT

Natural language supervision
Text-image retrieval
Webly supervised learning

# Try CLIP today!

https://github.com/openai/CLIP

- PyTorch implementation
- Colab notebook

# Thank You

Visit openai.com for more information.