



南京航空航天大学

Nanjing University of Aeronautics and Astronautics



模式分析与机器智能  
工业和信息化部重点实验室

MIT Key Laboratory of  
Pattern Analysis & Machine Intelligence

---

# Learning Transferable Visual Models From Natural Language Supervision

---

Alec Radford<sup>\*1</sup> Jong Wook Kim<sup>\*1</sup> Chris Hallacy<sup>1</sup> Aditya Ramesh<sup>1</sup> Gabriel Goh<sup>1</sup> Sandhini Agarwal<sup>1</sup>  
Girish Sastry<sup>1</sup> Amanda Askell<sup>1</sup> Pamela Mishkin<sup>1</sup> Jack Clark<sup>1</sup> Gretchen Krueger<sup>1</sup> Ilya Sutskever<sup>1</sup>

*OpenAI*

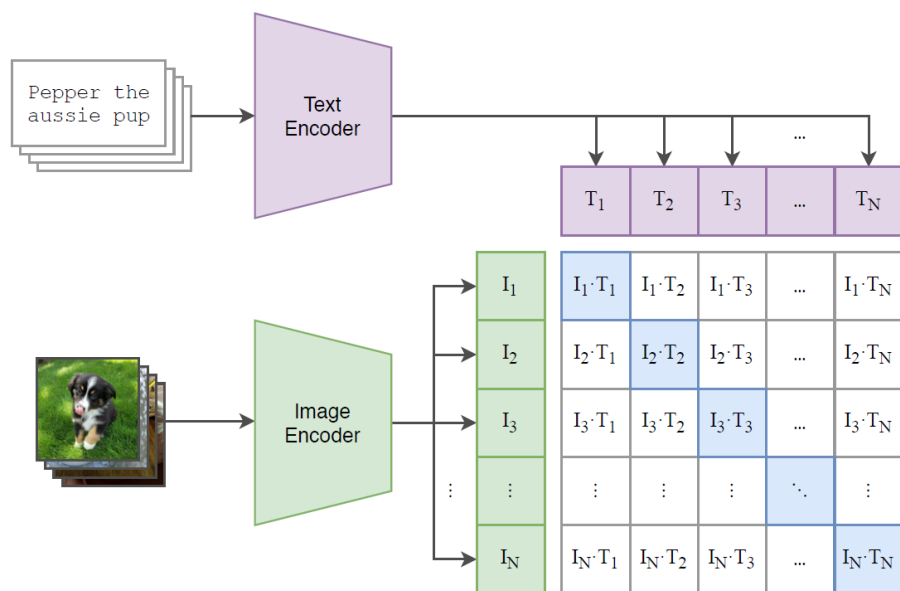
*ICML 2021*



# Framework

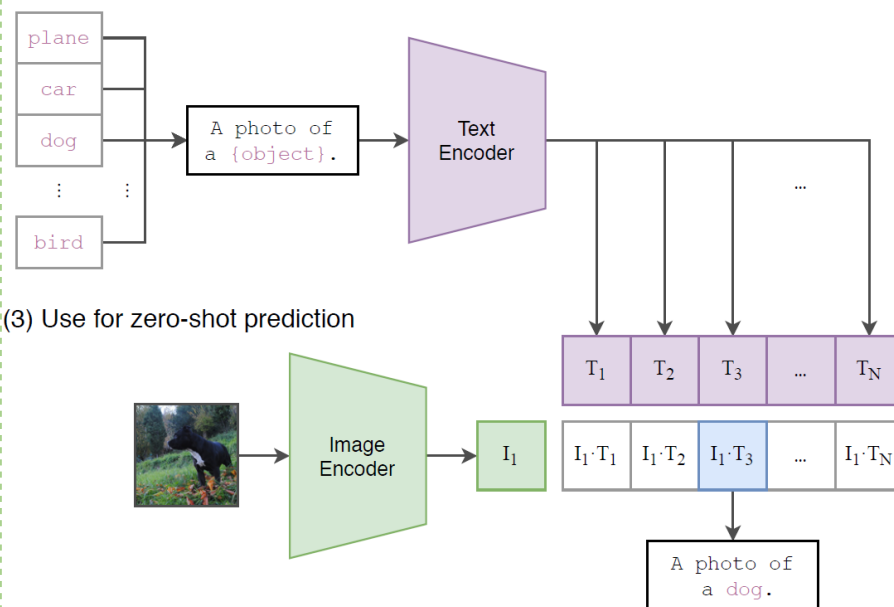
## Contrastive Language-Image Pre-training

(1) Contrastive pre-training

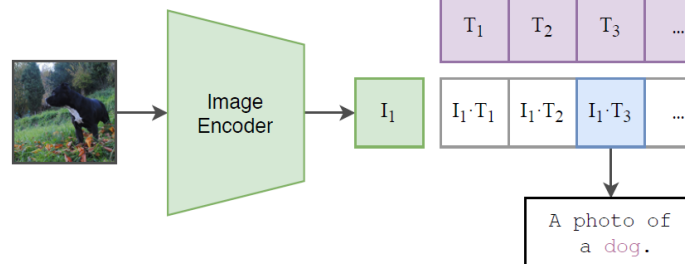


Contrastive pre-training

(2) Create dataset classifier from label text



(3) Use for zero-shot prediction

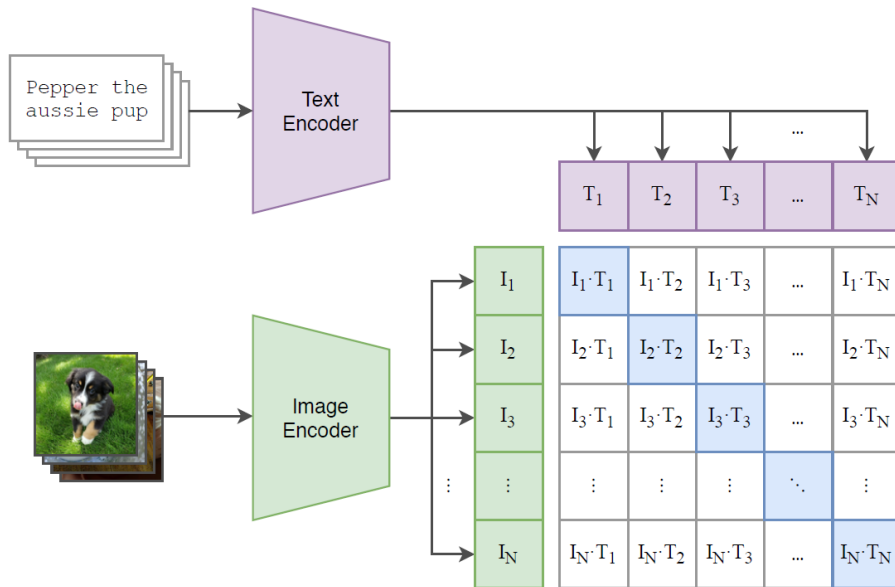


Zero-shot prediction

# Framework of CLIP

- Contrastive pre-training
  - Train an image encoder and a text encoder to predict the correct pairings of a batch of (image, text) training examples.

(1) Contrastive pre-training



Contrastive pre-training

```
# 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]

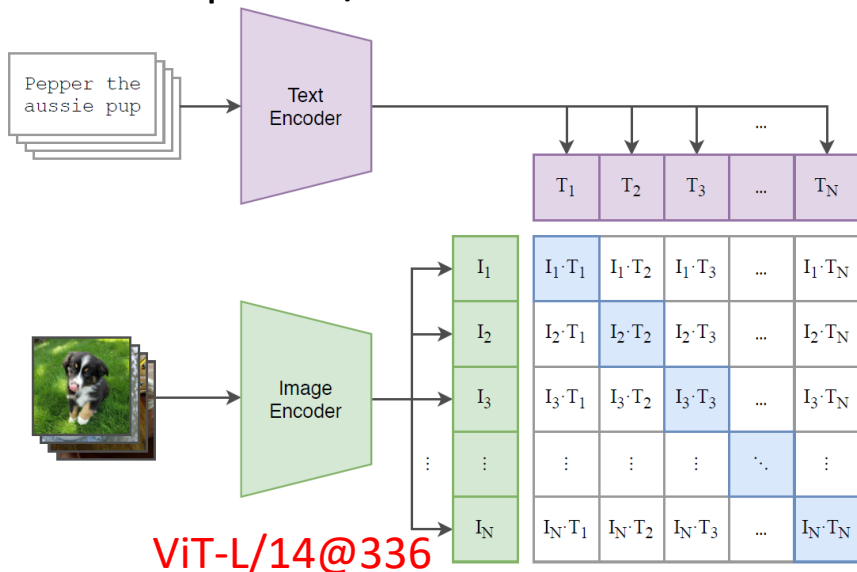
# joint multimodal embedding [n, d_e]
I_e = l2_normalize(np.dot(I_f, W_i), axis=1)
T_e = l2_normalize(np.dot(T_f, W_t), axis=1)

# 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
```

# Framework of CLIP

- Contrastive pre-training
  - Dataset: construct a new dataset of 400 million (image, text) pairs.
  - Text encoder: a text transformer (63M parameters)
  - Image encoder: 5 ResNets (ResNet-50, ResNet-101, RN50x4, RN50x16, RN50x64) and 3 vision transformer (ViT-B/32, ViT-B/16, ViT-L/14)
  - 32 epochs /32768 batch size



RN50X64: 18days (592 V100)

ViT-L/14: 12days (256 V100)

# Framework of CLIP

- Zero-shot prediction:
  - The learned text encoder synthesizes a zero-shot linear classifier by embedding the names or descriptions of the target dataset's classes.

(1) Create dataset classifier from label text

```

imagenet_templates = [
  'a bad photo of a {}.',
  'a photo of many {}.',
  'a sculpture of a {}.',
  'a photo of the hard to see {}.',
  'a low resolution photo of the {}.',
  'a rendering of a {}.',
  'graffiti of a {}.',
  'a bad photo of the {}.',
  'a cropped photo of the {}.',
  'a tattoo of a {}.',
  'the embroidered {}.',
  'a photo of a hard to see {}.',
  'a bright photo of a {}.',
  'a photo of a clean {}.',
  'a photo of a dirty {}.',
  'a dark photo of the {}.',
  'a drawing of a {}.',
  'a photo of my {}.',
]
  
```

↑ 3.5%

(2) Create dataset classifier from label text

plane, car, dog, ..., bird

↑ 1.3%

A photo of a {}.

Prompt template

Text Encoder

(3) Use for zero-shot prediction

Image Encoder

$I_1$

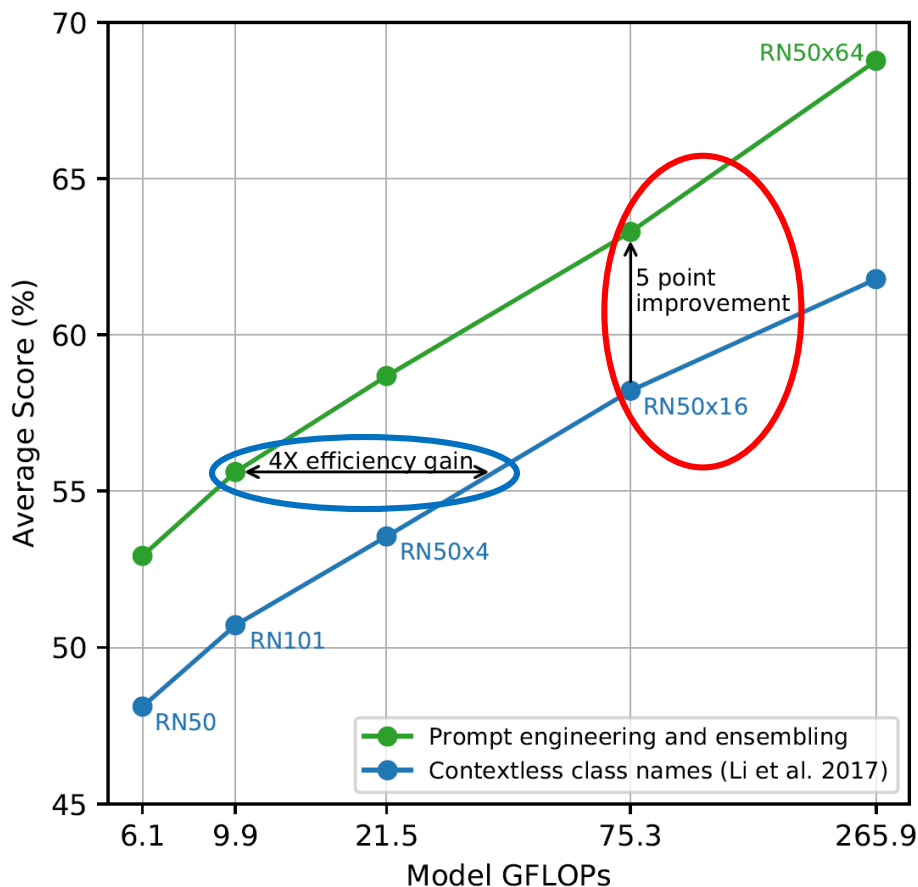
$T_1, T_2, T_3, \dots, T_N$

$I_1 \cdot T_1, I_1 \cdot T_2, I_1 \cdot T_3, \dots, I_1 \cdot T_N$

A photo of a dog.

Zero-shot prediction

# Prompt Engineering and Ensembling



**Figure 4. Prompt engineering and ensembling improve zero-shot performance.** Compared to the baseline of using contextless class names, prompt engineering and ensembling boost zero-shot classification performance by almost 5 points on average across 36 datasets. This improvement is similar to the gain from using 4 times more compute with the baseline zero-shot method but is “free” when amortized over many predictions.



# Experiments

# Initial Comparison on Visual N-Grams

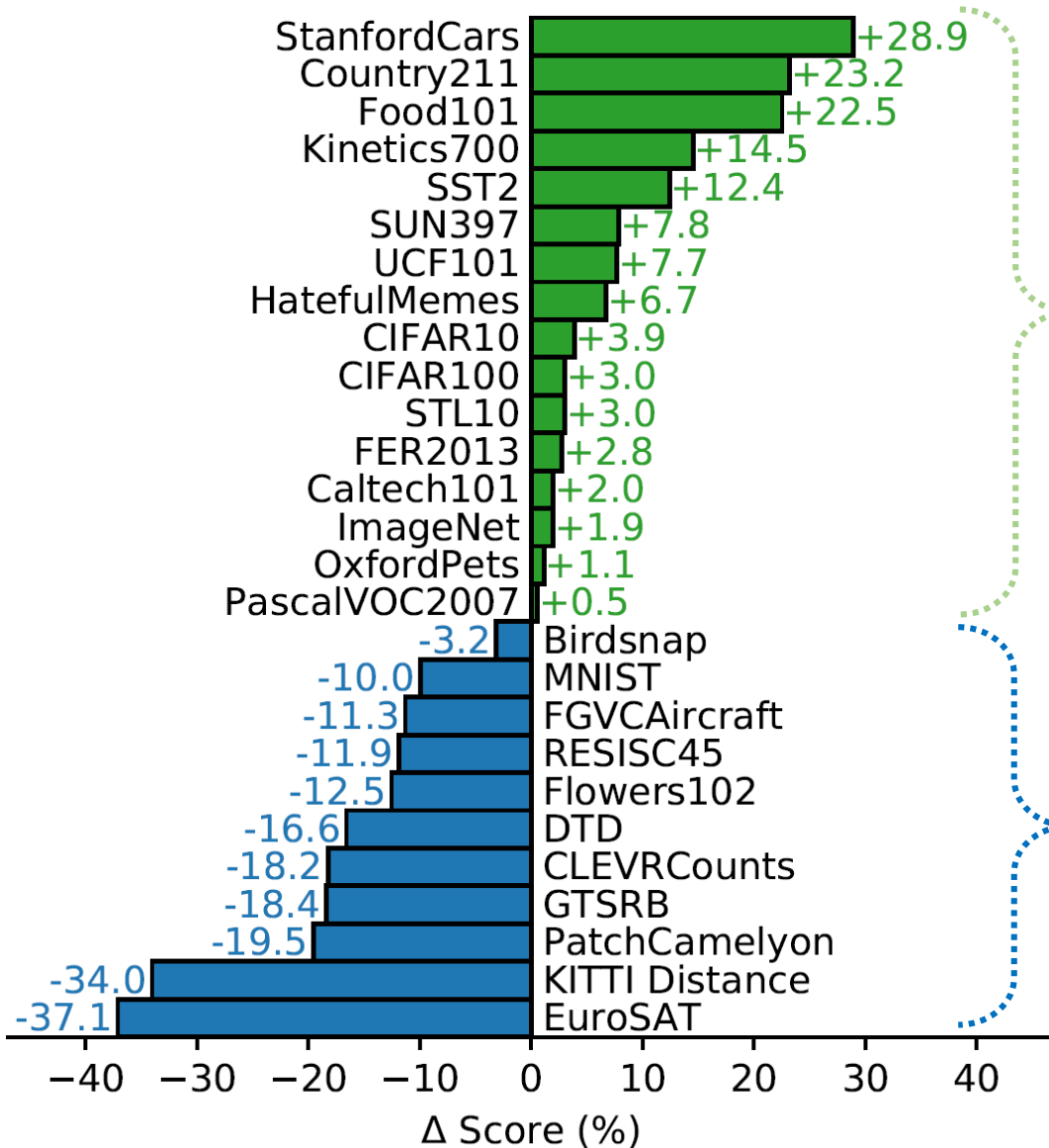


	aYahoo	ImageNet	SUN
Visual N-Grams	72.4	11.5	23.0
CLIP	<b>98.4</b>	<b>76.2</b>	<b>58.5</b>

Table 1. Comparing CLIP to prior zero-shot transfer image classification work. CLIP improves performance on all three datasets by a large amount. This improvement reflects many differences since the development of Visual N-Grams (Li et al., 2017).

- This is not a fair comparison
  - Dataset: 10x larger
  - Vision model: 100x more compute per prediction, use a transformer-based model.

# Zero-shot Learning

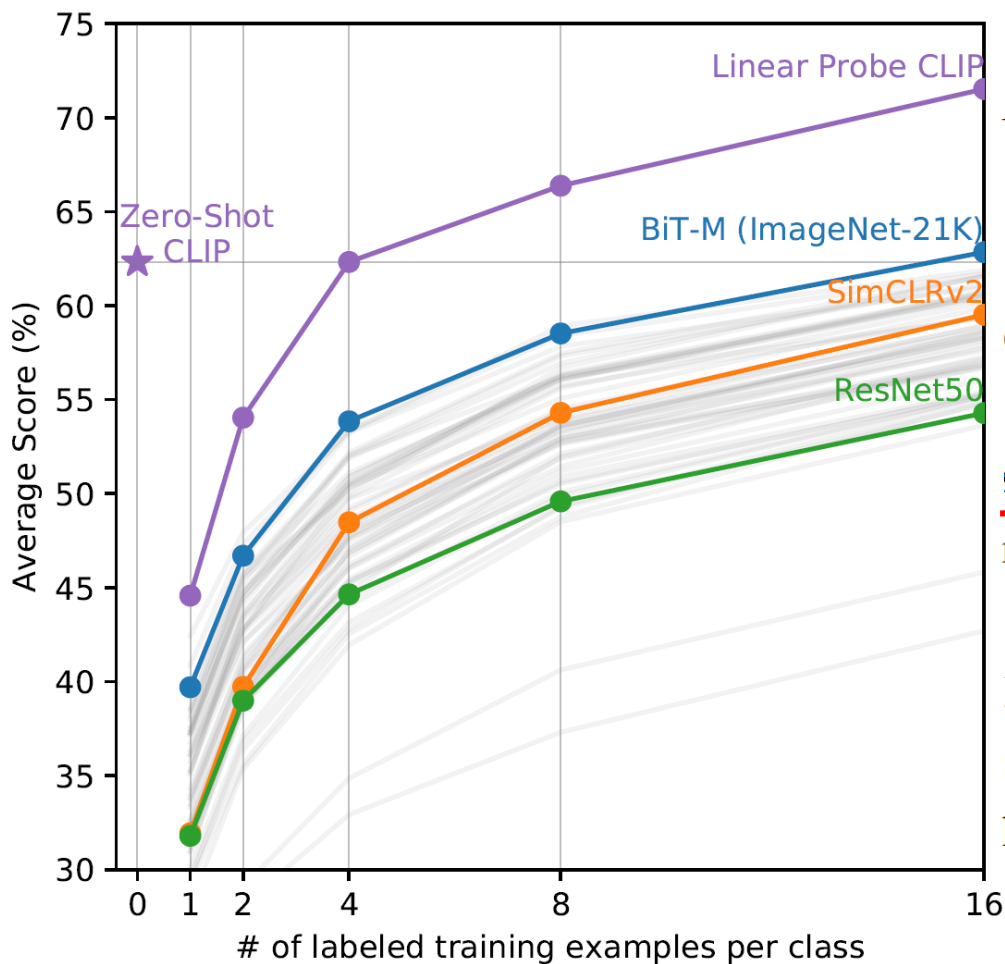


General tasks

Specialized, complex,  
abstract tasks

Zero-Shot CLIP vs. Linear Probe on ResNet50

# Few-shot Learning

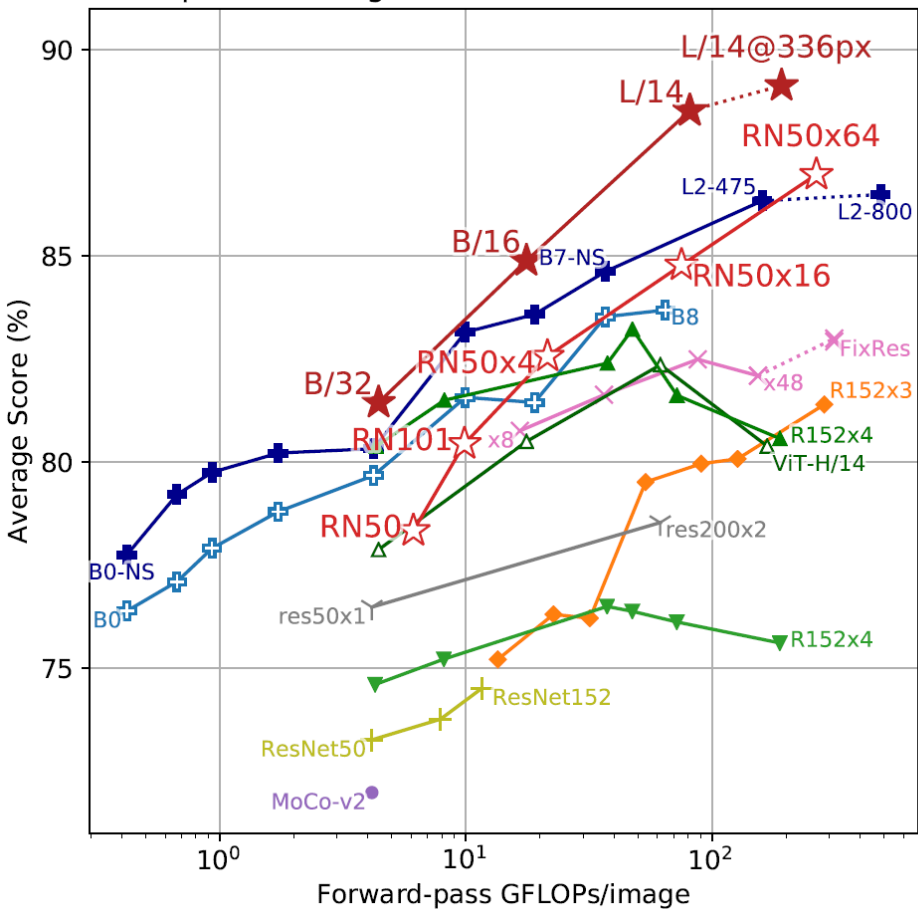


*Figure 6. Zero-shot CLIP outperforms few-shot linear probes. Zero-shot CLIP matches the average performance of a 4-shot linear classifier trained on the same feature space and nearly matches the best results of a 16-shot linear classifier across publicly available models. For both BiT-M and SimCLRv2, the best performing model is highlighted. Light gray lines are other models in the eval suite. The 20 datasets with at least 16 examples per class were used in this analysis.*

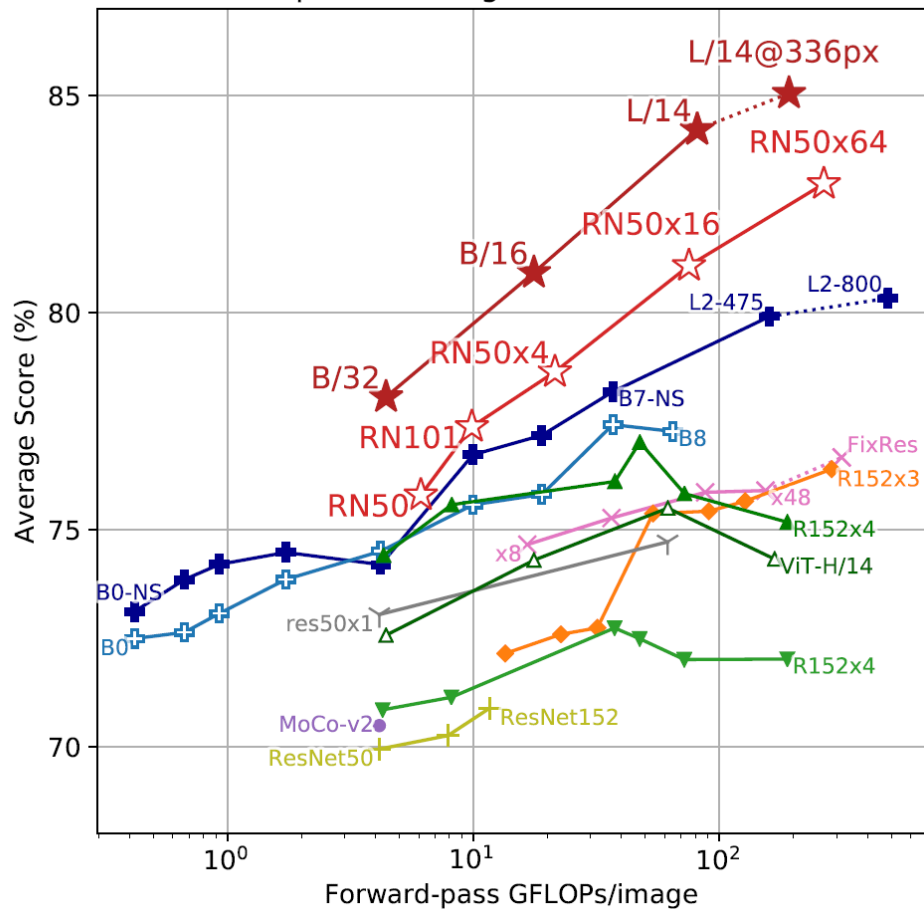


# Representation Learning

Linear probe average over Kornblith et al.'s 12 datasets



Linear probe average over all 27 datasets



- ★ CLIP-ViT
- ✕ Instagram-pretrained
- ▲ ViT (ImageNet-21k)
- ☆ CLIP-ResNet
- ◆ SimCLRv2
- ▲ BiT-M
- EfficientNet-NoisyStudent
- BYOL
- ▼ BiT-S
- + EfficientNet
- MoCo
- + ResNet

# Robust to distribution shift



	Dataset Examples	ImageNet ResNet101	Zero-Shot CLIP	$\Delta$ 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%

Figure 13. **Zero-shot CLIP is much more robust to distribution shift than standard ImageNet models.**

(Right) Visualizing distribution shift for bananas, a class shared across 5 of the 7 natural distribution shift datasets. The performance of the best zero-shot CLIP model, ViT-L/14@336px, is compared with a model that has the same performance on the ImageNet validation set, ResNet-101.



# Limitations

# Framework of CLIP



The performance of zero-shot CLIP is on average competitive with the simple supervised baseline, not SOTA.

The performance of zero-shot CLIP is still quite weak on several kinds of tasks, such as fine-grained classification.

The zero-shot CLIP generalizes poorly to data that is truly out-of-distribution for it, such as MNIST.

CLIP is still limited to choosing from only those concepts in a given zero-shot classifier.



# Motivation

# Motivation



- Pre-training of CV
  - Supervised Learning
  - Self-supervised Learning: MOCO/SimCLR, MAE/BeiT

The pre-training model still needs supervised fine-tuning when transferring to downstream tasks, and cannot achieve zero-shot.

- Pre-training of NLP
  - Task-agnostic objectives: autoregressive and masked language modeling.
  - Zero-shot transfer to downstream datasets.

Could scalable pre-training methods which learn directly from web text result in a similar breakthrough in computer vision ?

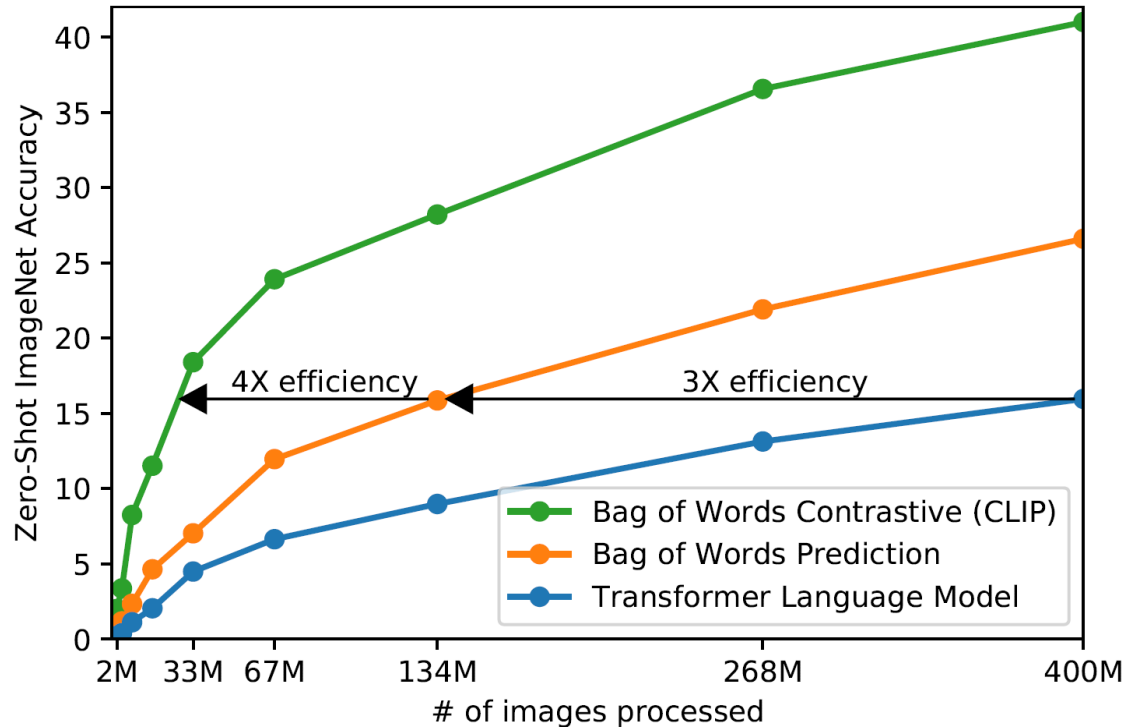
# Related work

- Natural language supervision
  - Visual N-grams: zero-shot transfer [Li et al., ICCV'17]
  - VirTex: transform-based language modeling [Desai & Johnson et al., ICCV'21]
  - ICMLM: masked-language modeling [Bulent Sariyildiz et al., ECCV'20]
  - ConVIRT: contrastive objectives [Zhang et al., arXiv'20]
- Natural language weak supervision
  - BiT: [Kolesnikov et al., ECCV'20]
  - ViT: [Dosovitskiy et al., ICCV'21]

JFT-300M/18291 classes

A crucial difference between these weakly supervised models and recent explorations of learning image representations directly from natural language is scale.

# Contrastive Pre-training



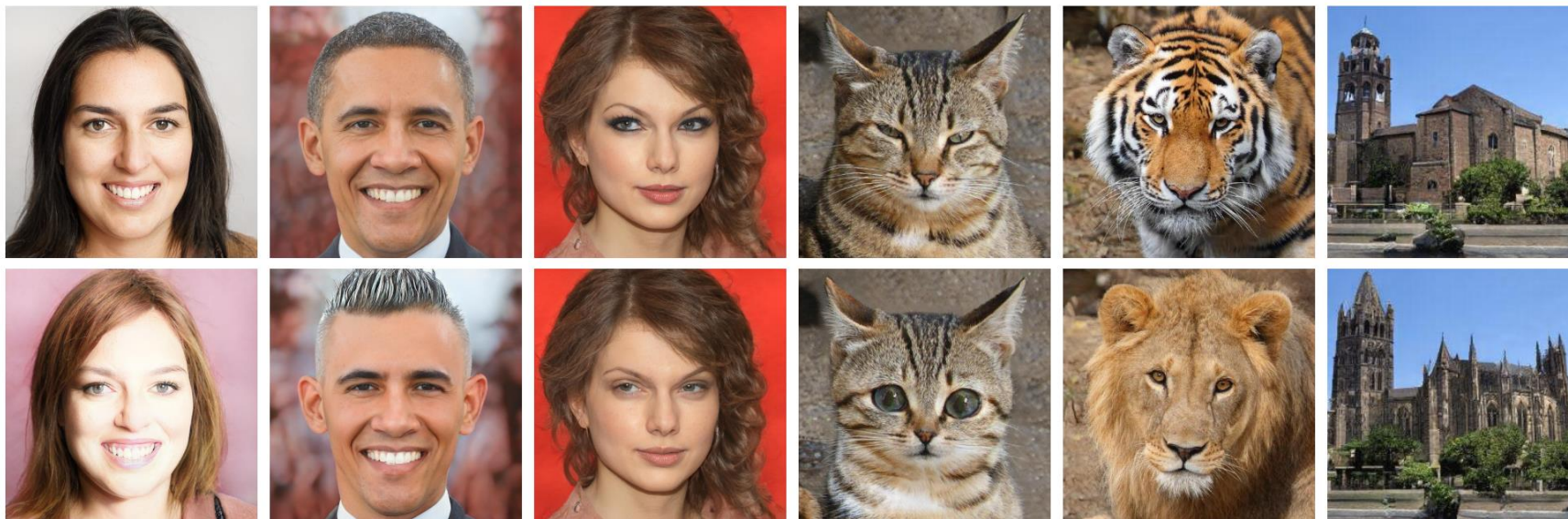
*Figure 2. CLIP is much more efficient at zero-shot transfer than our image caption baseline. Although highly expressive, we found that transformer-based language models are relatively weak at zero-shot ImageNet classification. Here, we see that it learns 3x slower than a baseline which predicts a bag-of-words (BoW) encoding of the text (Joulin et al., 2016). Swapping the prediction objective for the contrastive objective of CLIP further improves efficiency another 4x.*



**Application**

## StyleCLIP: Text-Driven Manipulation of StyleGAN Imagery

Or Patashnik<sup>†\*</sup>   Zongze Wu<sup>‡\*</sup>   Eli Shechtman<sup>§</sup>   Daniel Cohen-Or<sup>†</sup>   Dani Lischinski<sup>‡</sup>  
<sup>‡</sup>Hebrew University of Jerusalem   <sup>†</sup>Tel-Aviv University   <sup>§</sup>Adobe Research



“Emma Stone”

“Mohawk hairstyle”

“Without makeup”

“Cute cat”

“Lion”

“Gothic church”

Figure 1. Examples of text-driven manipulations using StyleCLIP. Top row: input images; Bottom row: our manipulated results. The text prompt used to drive each manipulation appears under each column.



**a man**



**a handsome man**



**a blonde man**

## CLIPDraw: Exploring Text-to-Drawing Synthesis through Language-Image Encoders

Kevin Frans<sup>1,2</sup>, L. B. Soros<sup>1</sup> and Olaf Witkowski<sup>1,3,4</sup>



“A drawing of a cat”.



“Horse eating a cupcake”.



“A 3D rendering of a temple”.



“Family vacation to Walt Disney World”.



“Self”.

**Various drawings synthesized by CLIPDraw**, along with the corresponding description prompts used. CLIPDraw synthesizes images from text by performing gradient descent over a set of RGBA Bézier curves, with the goal of minimizing cosine distance between the CLIP encodings of generated images and description prompts. CLIPDraw does not require learning a new model, and can generally synthesize images within a minute on a typical GPU.

*Frans et al. CLIPDraw: Exploring Text-to-Drawing Synthesis through Language-Image Encoders. arXiv, 2021*

## Contrastive Language-Image Forensic Search



A truck with the text "odwalla"



A white BMW car



A bicyclist with a blue shirt



A blue SMART car



**Thanks!**