

# Learning Transferable Visual Models From Natural Language Supervision

## - CLIP(Contrastive Language-Image Pre-training)

### Abstract

- State-of-the-art computer vision systems are trained to predict a fixed set of predetermined object categories.
- This restricted form of supervision limits their generality and usability since additional labeled data is needed to specify any other visual concept.
- Learning directly from raw text about images is a promising alternative which leverages a much broader source of supervision.
- We demonstrate that the simple pre-training task of predicting which caption goes with which image is an efficient and scalable way to learn SOTA image representations from scratch on a dataset of 400 million (image, text) pairs collected from the internet.
- After pre-training, natural language is used to reference learned visual concepts (or describe new ones) enabling zero-shot transfer of the model to downstream tasks.
- We study the performance of this approach by benchmarking on over 30 different existing computer vision datasets, spanning tasks such as OCR, action recognition in videos, geo-localization, and many types of fine-grained object classification.
- The model transfers non-trivially to most tasks and is often competitive with a fully supervised baseline without the need for any dataset specific training.
- For instance, we match the accuracy of the original ResNet-50 on ImageNet zero-shot without needing to use any of the 1.28 million training examples it was trained on.
- We release our code and pre-trained model weights at <https://github.com/OpenAI/CLIP> (<https://github.com/OpenAI/CLIP>).

### Motivation

(의역 많음)

- NLP 분야에서는 Pre-training method로 인해 많이 발전했다.
- Web에서의 많은(web scale) NLP data로 인해 그러할 것이다.
- 다른 field에서는 그렇지 못하다(vision에서는 여전히 ImageNet과 같은 crowd-labeled datasets으로 pretrain model을 기준으로 한다.)
- 이러한 Web scale을 통해 computer vision에서의 발전을 할 수 있을 것이다.

### Approach

#### 1. Natural Language Supervision

- Learning from natural language has several potential strengths over other training methods.

#### 2. Dataset

## 1. Previous text-image dataset

- MS-COCO and Visual Genome: high quality crowd-labeled datasets. 100,000
- TFCC100M: 100 million photos but sparse and varying quality. -> 가공후 6~15 million.

## 2. New dataset (WIT; WebImageText)

- 400 million image text pairs
- class balance version : 20000 pairs per query

## 3. Training Method

### 1. initial approach

- CNN + Transformer like VirTex

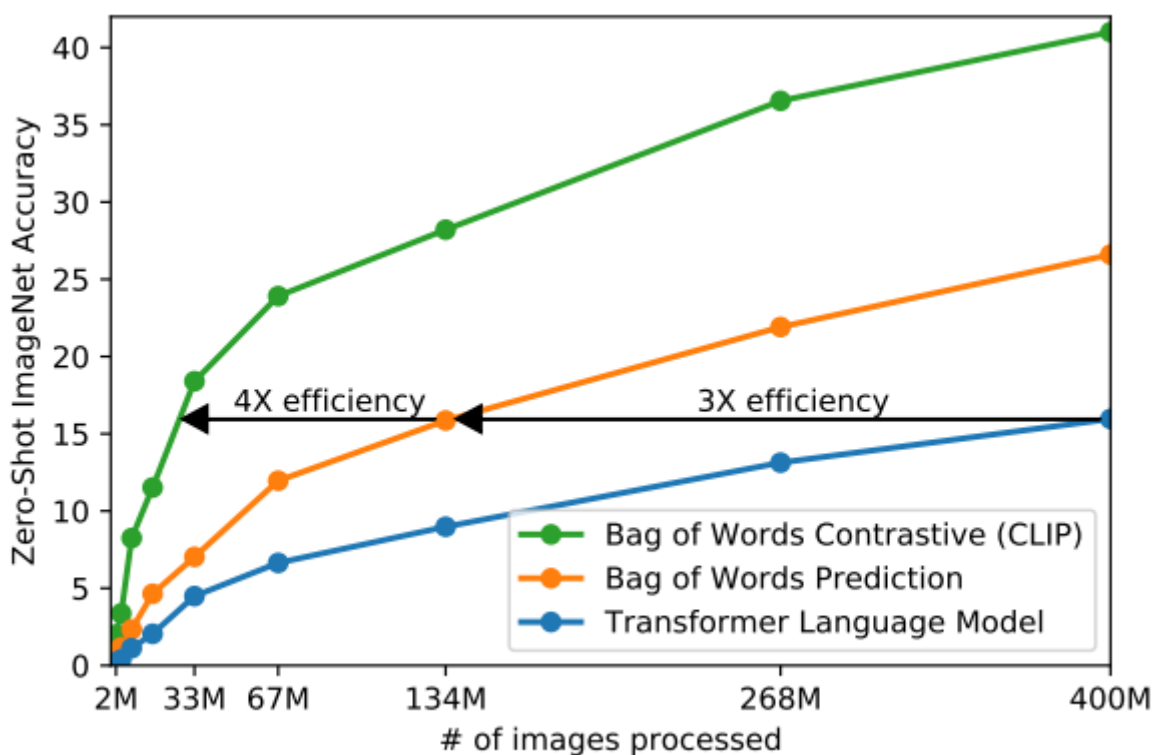
### 2. Bag of Words Prediction

- Using Bag of Words instead of pull 'exact' text prediction

### 3. Bag of words contrastive(CLIP)

- contrastive learning
- predicting only which text as a whole is paired with which image and not the exact words of that text.

## Result



Constrastive 방법이 매우 효율적이다.

**VirText**

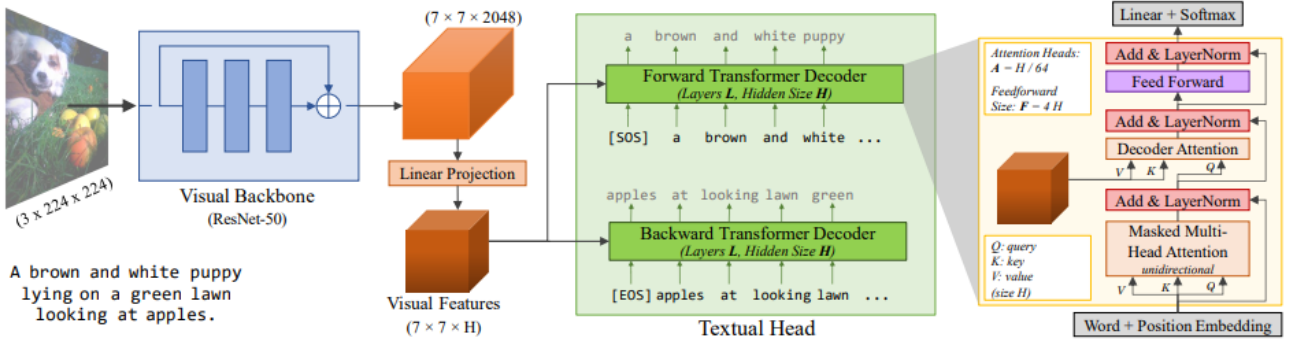
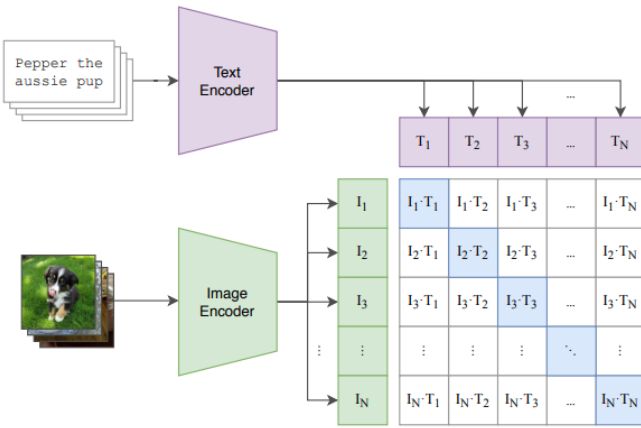


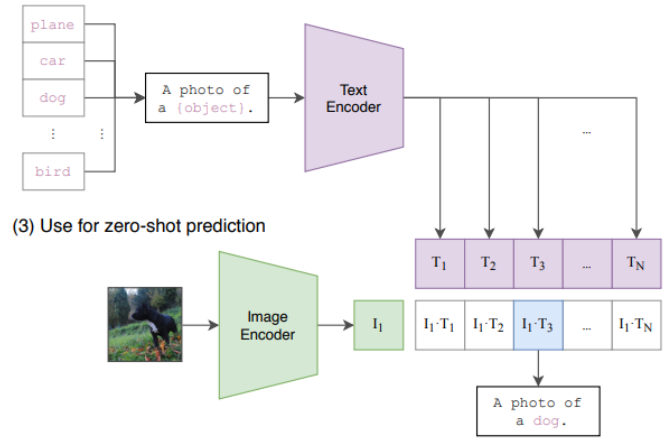
Figure 3: **VirTex pretraining setup:** Our model consists of a *visual backbone* (ResNet-50), and a *textual head* (two uni-directional Transformers). The visual backbone extracts image features, and textual head predicts captions via bidirectional language modeling (*bicaptioning*). The Transformers perform masked multiheaded self-attention over caption features, and multiheaded attention over image features. Our model is trained end-to-end from scratch. After pretraining, the visual backbone is transferred to downstream visual recognition tasks.

## CLIP

### (1) Contrastive pre-training



### (2) Create dataset classifier from label text



### (3) Use for zero-shot prediction

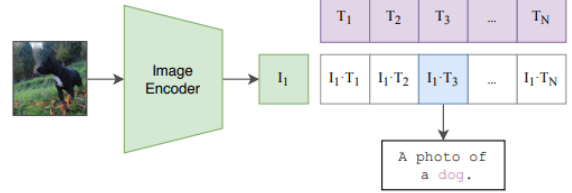


Figure 1. Summary of our approach. While standard image models jointly train an image feature extractor and a linear classifier to predict some label, CLIP jointly trains an image encoder and a text encoder to predict the correct pairings of a batch of (image, text) training examples. At test time the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

## Experiments

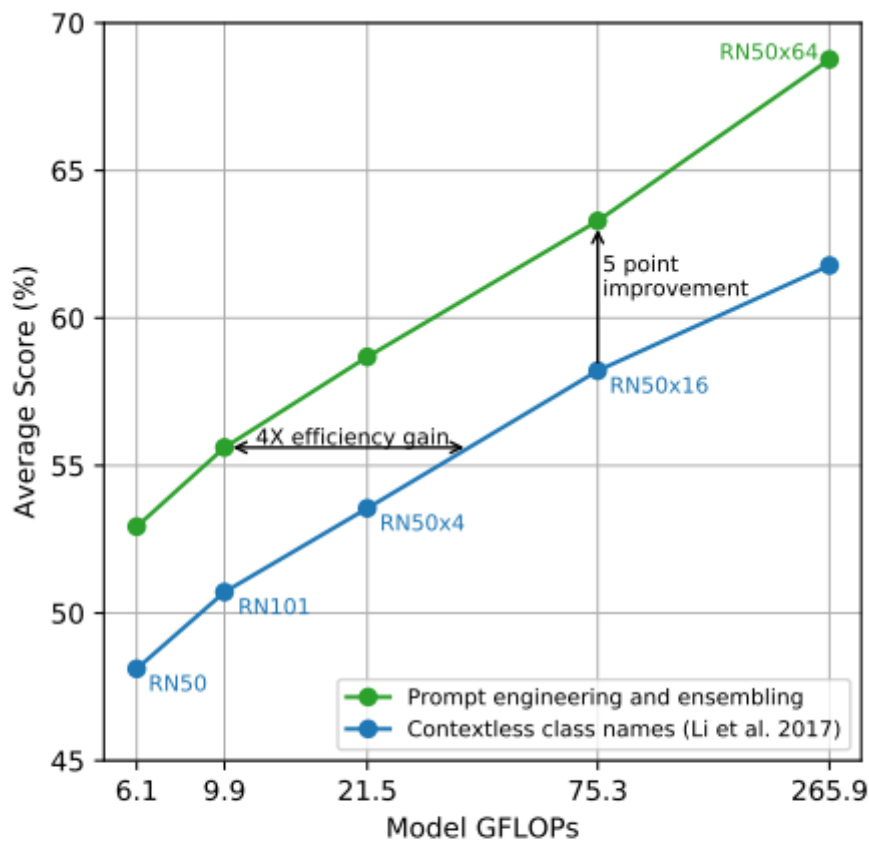
### 1. ZeroShot transfer

#### Compare with Visual N-Grams(2017)

	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>

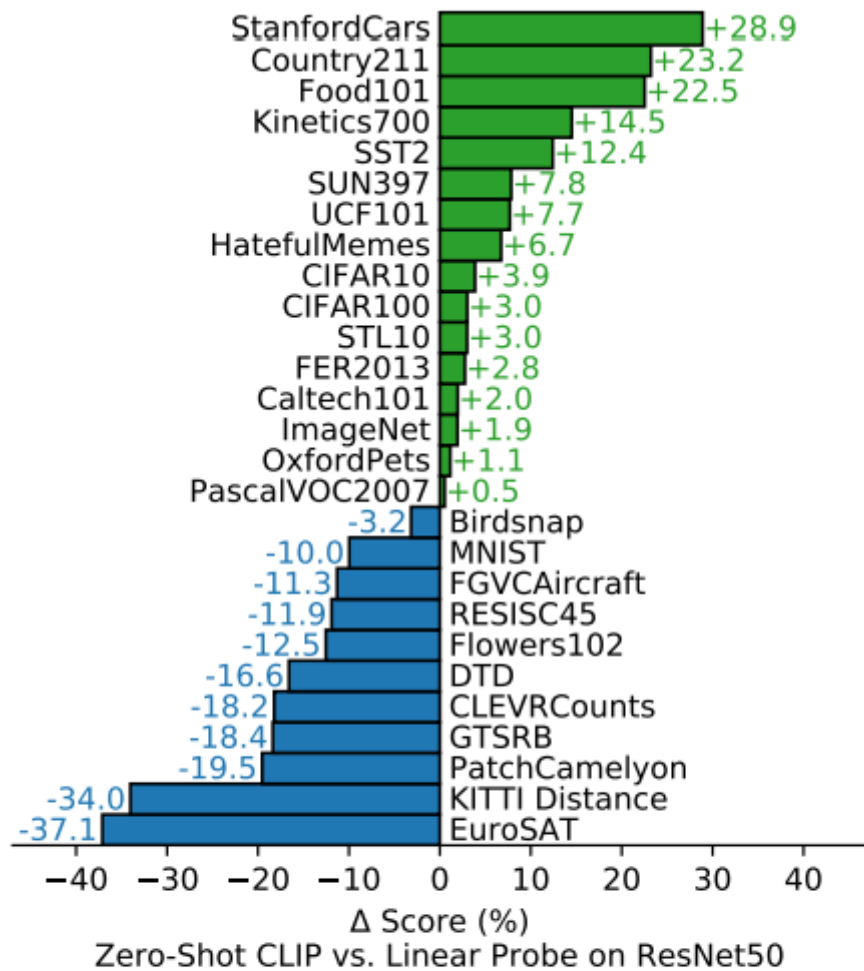
*Table 1.* Comparing CLIP to prior zero-shot transfer image classification results. CLIP improves performance on all three datasets by a large amount. This improvement reflects many differences in the 4 years since the development of Visual N-Grams (Li et al., 2017).

### Ensembling improve



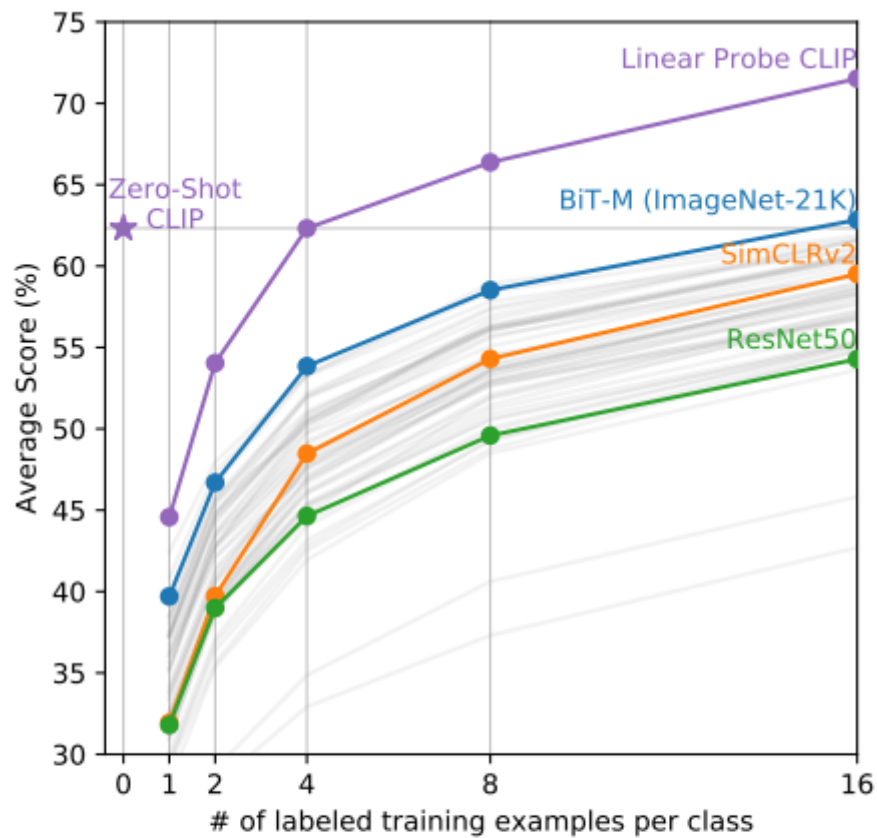
**Figure 4. Prompt engineering and ensembling improve zero-shot performance.** Compared to the baseline of using contextless class names, prompt engineering and ensembling boost zero-shot classification performance by almost 5 points on average across 36 datasets. This improvement is similar to the gain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

vs ResNet50

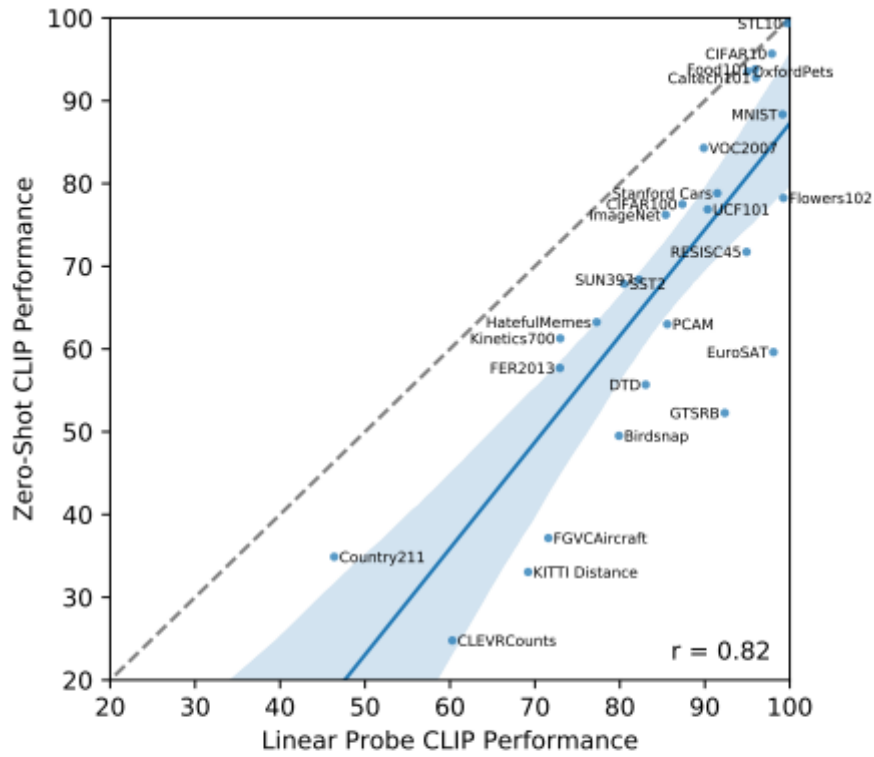


**Figure 5. Zero-shot CLIP is competitive with a fully supervised baseline.** Across a 27 dataset eval suite, a zero-shot CLIP classifier outperforms a fully supervised linear classifier fitted on ResNet-50 features on 16 datasets, including ImageNet.

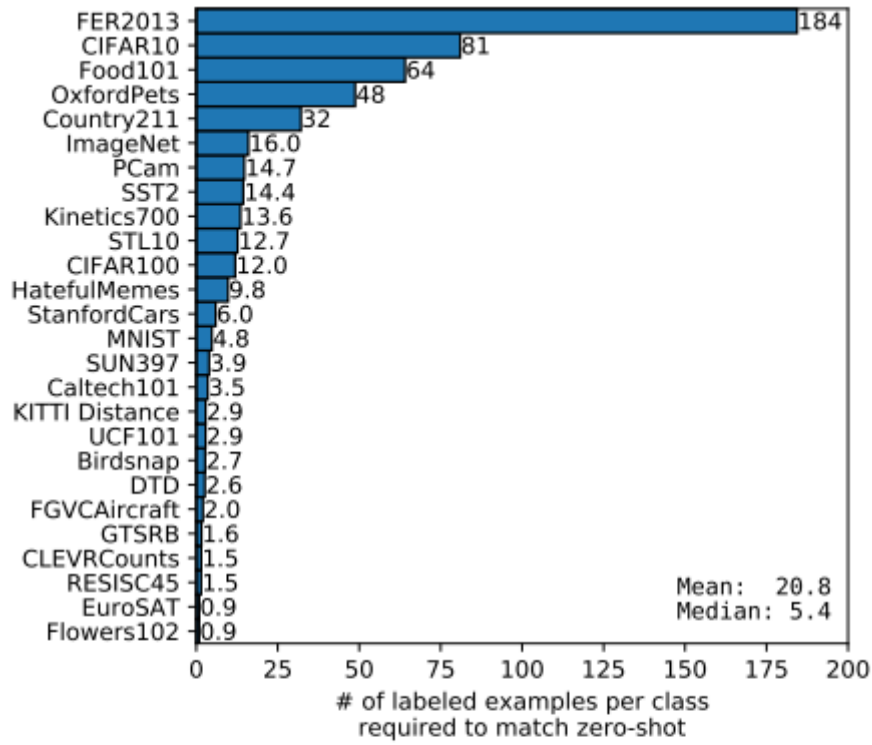
### ***Zeroshot vs Few-shot***



**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.

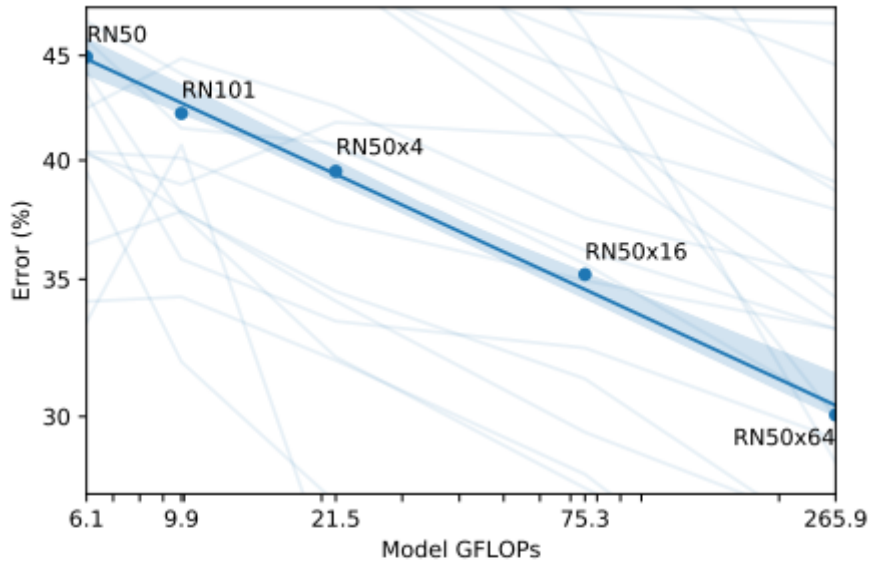


**Figure 8. Zero-shot performance is correlated with linear probe performance but still mostly sub-optimal.** Comparing zero-shot and linear probe performance across datasets shows a strong correlation with zero-shot performance mostly shifted 10 to 25 points lower. On only 5 datasets does zero-shot performance approach linear probe performance ( $\leq 3$  point difference).



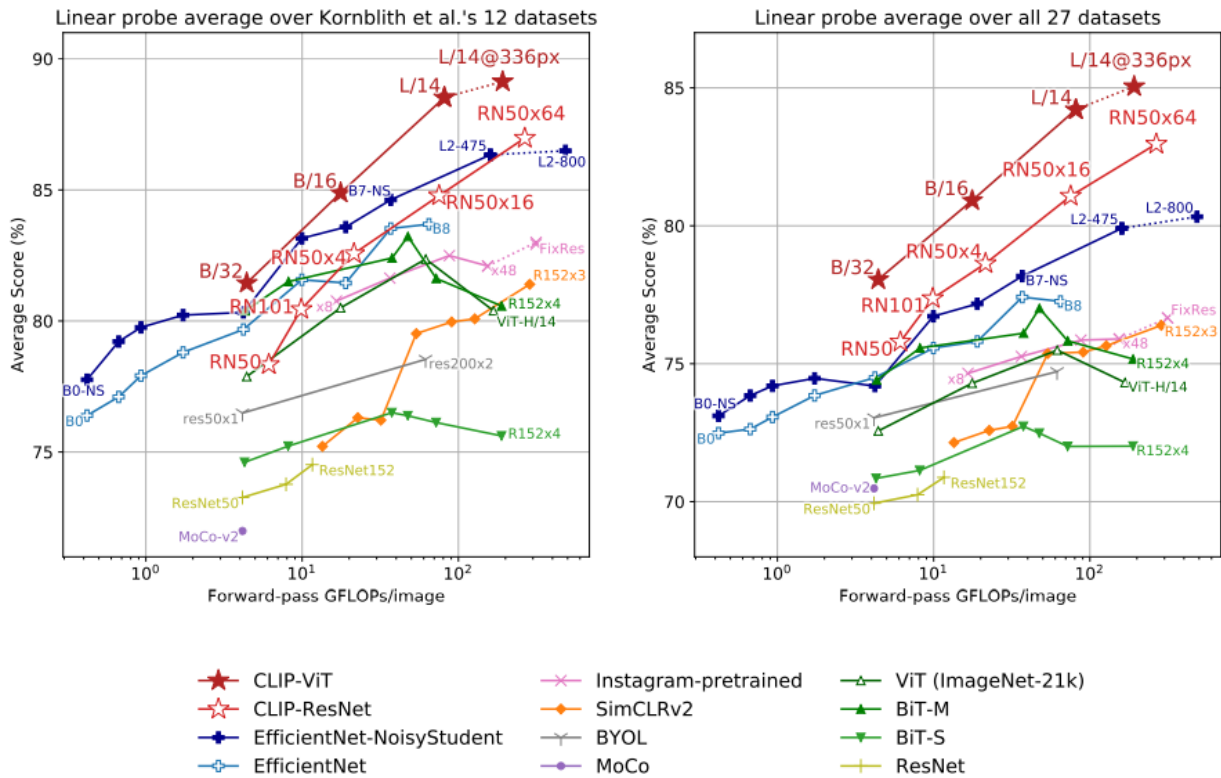
**Figure 7. The data efficiency of zero-shot transfer varies widely.** Calculating the number of labeled examples per class a linear classifier on the same CLIP feature space requires to match the performance of the zero-shot classifier contextualizes the effectiveness of zero-shot transfer. Values are estimated based on log-linear interpolation of 1, 2, 4, 8, 16-shot and fully supervised results. Performance varies widely from still underperforming a one-shot classifier on two datasets to matching an estimated 184 labeled examples per class.



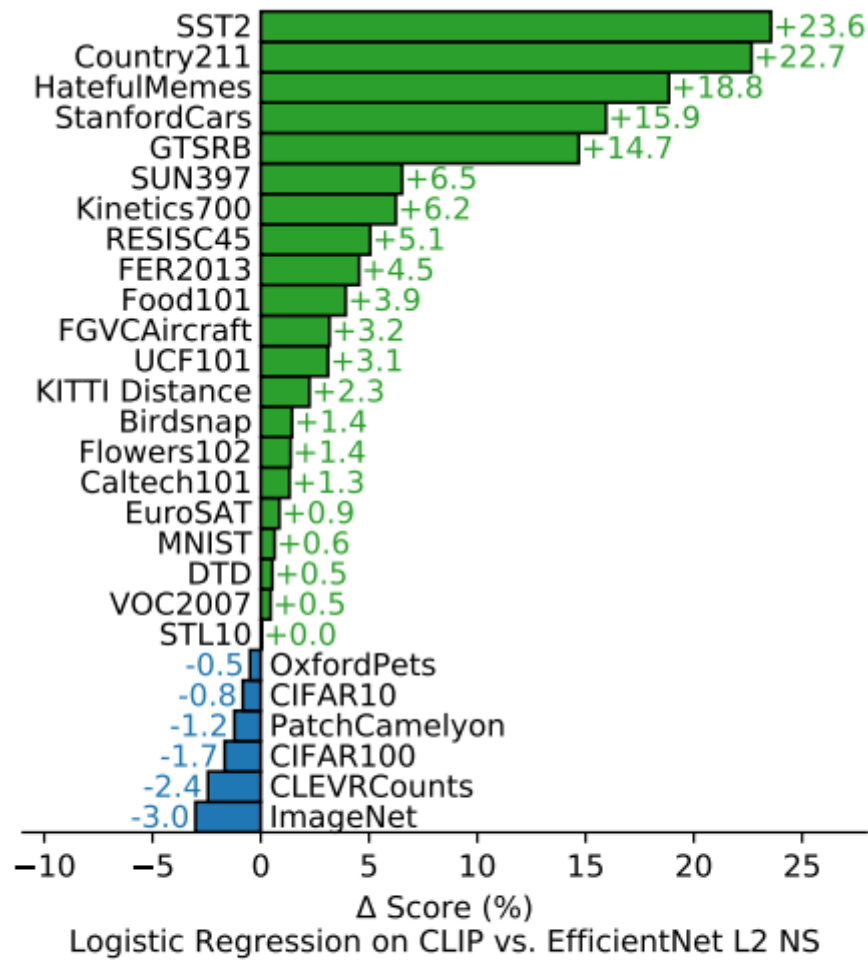


**Figure 9. Zero-shot CLIP performance scales smoothly as a function of model compute.** Across 39 evals on 36 different datasets, average zero-shot error is well modeled by a log-log linear trend across a 44x range of compute spanning 5 different CLIP models. Lightly shaded lines are performance on individual evals, showing that performance is much more varied despite the smooth overall trend.

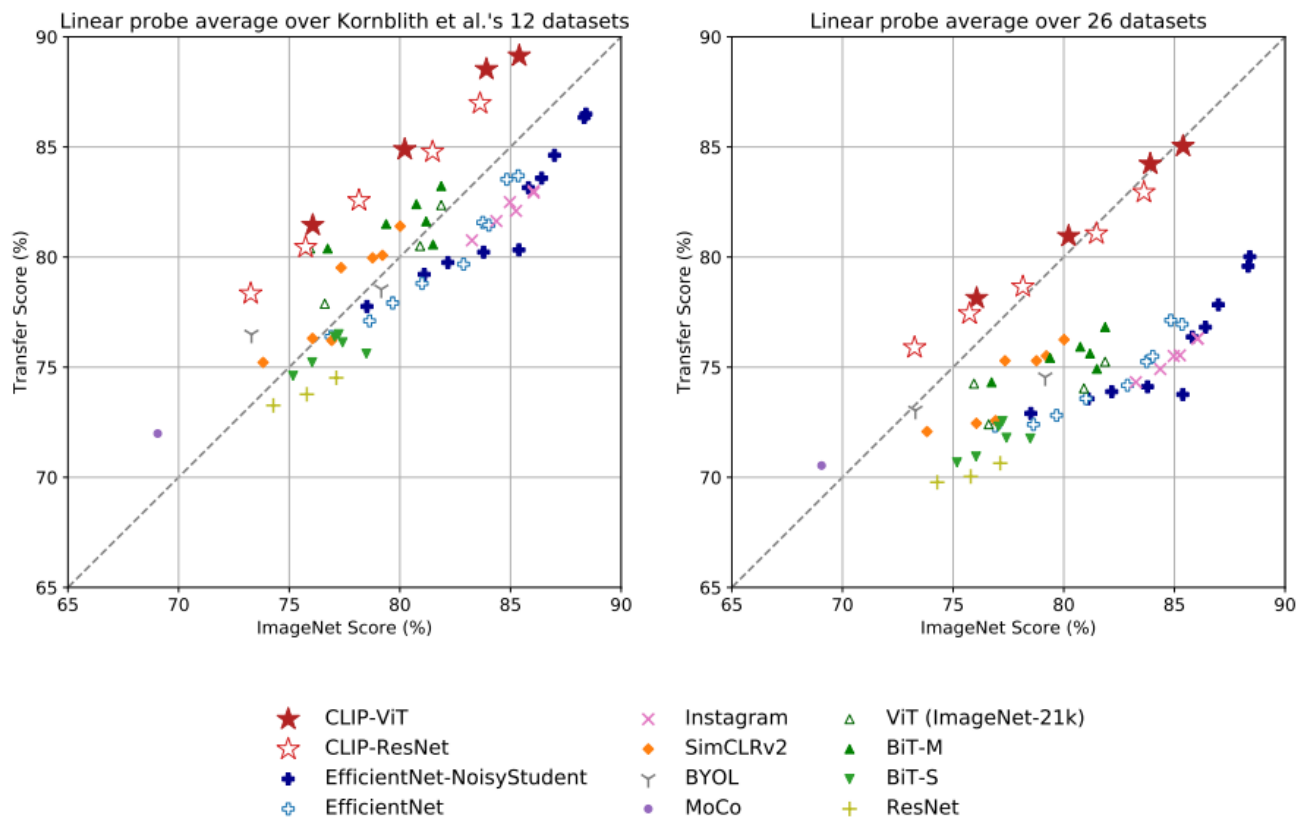
## 2. Representation Learning(Linear Probe)



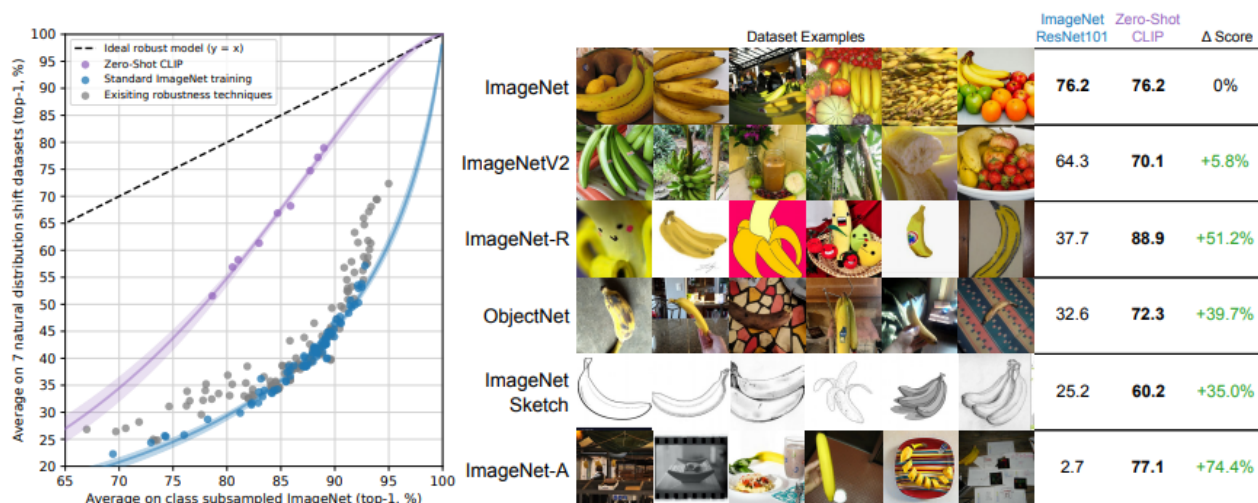
**Figure 10. Linear probe performance of CLIP models in comparison with state-of-the-art computer vision models**, including EfficientNet (Tan & Le, 2019; Xie et al., 2020), MoCo (Chen et al., 2020d), Instagram-pretrained ResNeXt models (Mahajan et al., 2018; Touvron et al., 2019), BiT (Kolesnikov et al., 2019), ViT (Dosovitskiy et al., 2020), SimCLRv2 (Chen et al., 2020c), BYOL (Grill et al., 2020), and the original ResNet models (He et al., 2016b). (Left) Scores are averaged over 12 datasets studied by Kornblith et al. (2019). (Right) Scores are averaged over 27 datasets that contain a wider variety of distributions. Dotted lines indicate models fine-tuned or evaluated on images at a higher-resolution than pre-training. See Table 10 for individual scores and Figure 20 for plots for each dataset.



**Figure 11. CLIP’s features outperform the features of the best ImageNet model on a wide variety of datasets.** Fitting a linear classifier on CLIP’s features outperforms using the Noisy Student EfficientNet-L2 on 21 out of 27 datasets.



**Figure 12. CLIP’s features are more robust to task shift when compared to models pre-trained on ImageNet.** For both dataset splits, the transfer scores of linear probes trained on the representations of CLIP models are higher than other models with similar ImageNet performance. This suggests that the representations of models trained on ImageNet are somewhat overfit to their task.



**Figure 13. Zero-shot CLIP is much more robust to distribution shift than standard ImageNet models.** (Left) An ideal robust model (dashed line) performs equally well on the ImageNet distribution and on other natural image distributions. Zero-shot CLIP models shrink this “robustness gap” by up to 75%. Linear fits on logit transformed values are shown with bootstrap estimated 95% confidence intervals. (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.

In [ ]:

