

SEMI-SUPERVISED LEARNING IN VISION

Already covered...

We have already covered (in the OCR case):

- pseudo-labels
- active learning
- data augmentation and prior knowledge of invariances
- EM algorithms

Pre-Transformer Approaches

Prior attempts to carry over principles for language models and HMMs to the image domain:

- linearize images and apply HMM or LSTM models
- apply VQ to patches and apply syntactic models to the resulting "visual words"
- corrupt images with noise and training a network to restore them
- predict color images from grayscale
- mask parts of images and predict the masked parts (like BERT)
- determine the spatial relations between patches (like entailment)

All of these yield deep learning architectures that are potentially useful for transfer learning, but never beat supervised models.

Masked Predictions



(a) Input context

(b) Human artist

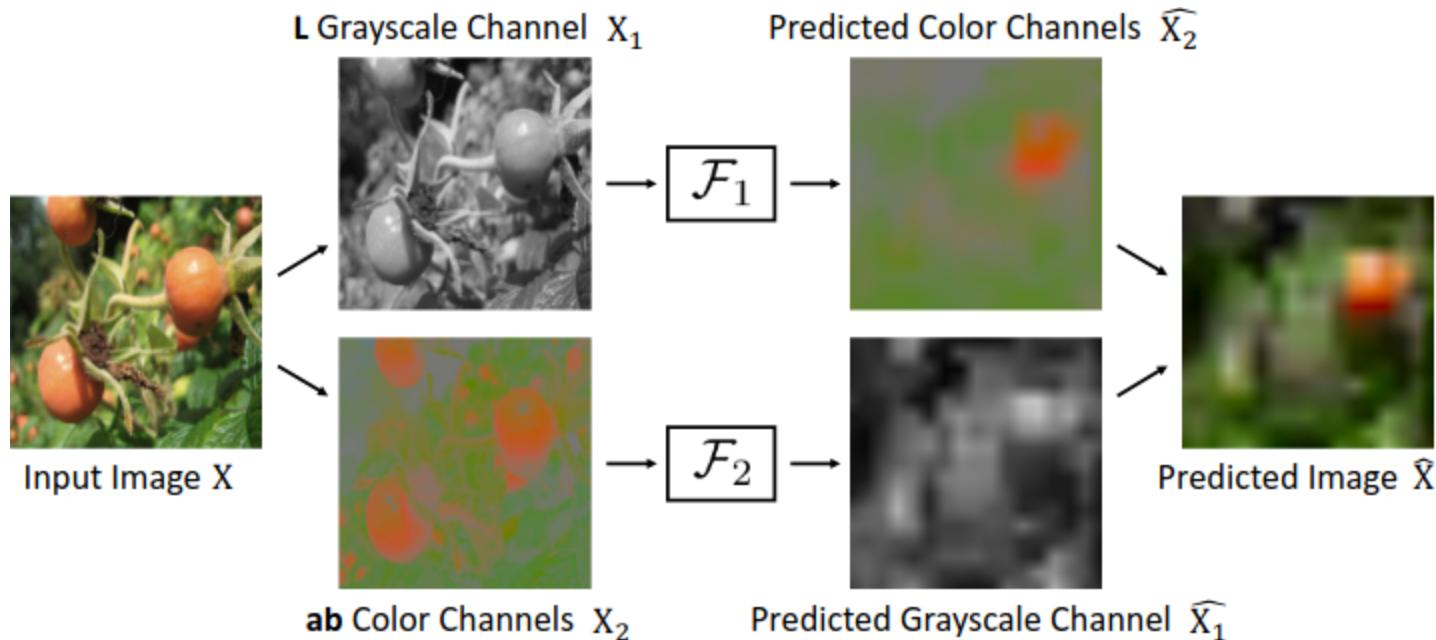


(c) Context Encoder
(L_2 loss)

(d) Context Encoder
(L_2 + Adversarial loss)

Pathak et al. 2016

Split-Brain Autoencoder



Zhang et al. 2016

Context Encoding

Example:



Question 1:



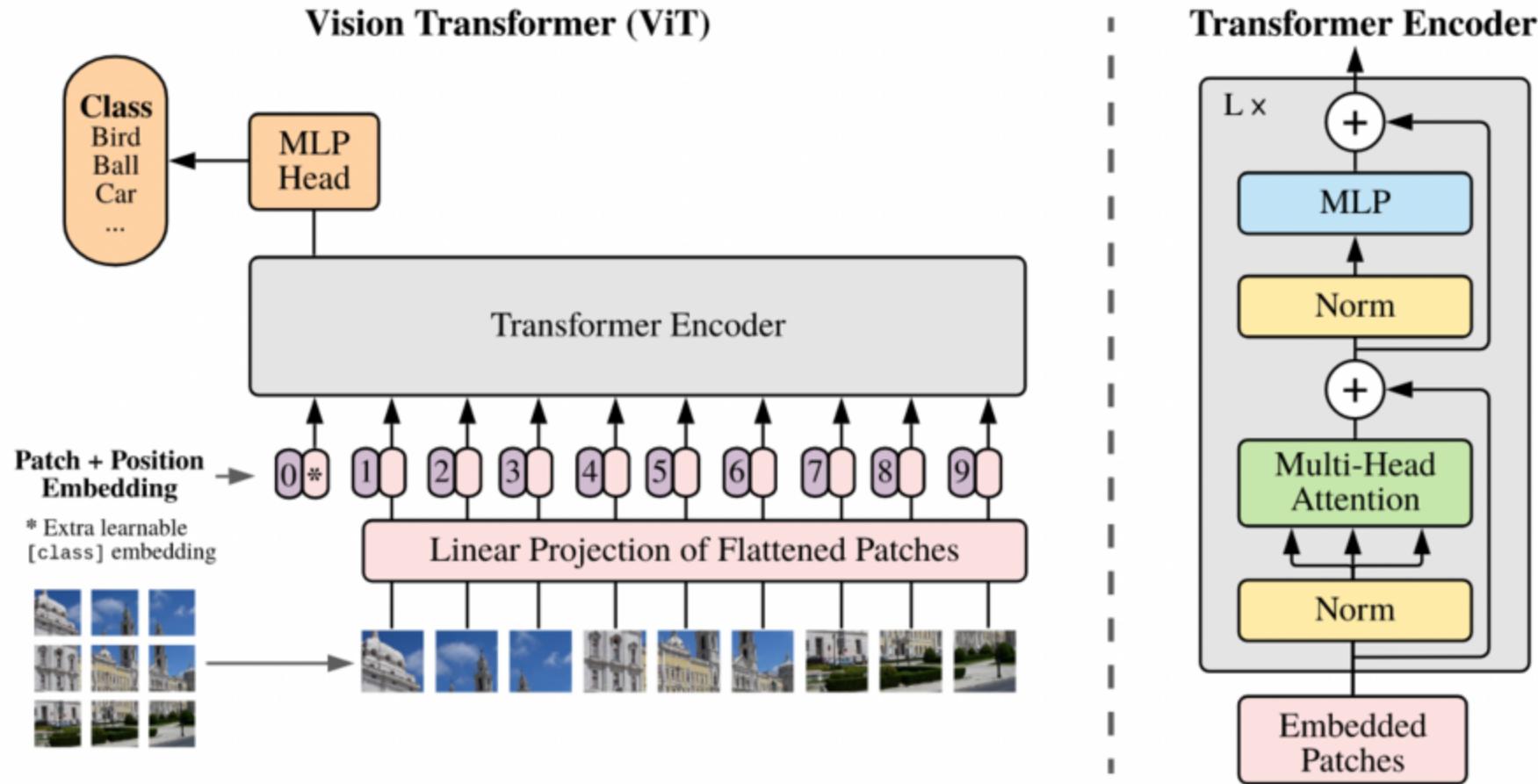
Question 2:



Doersch et al. 2016

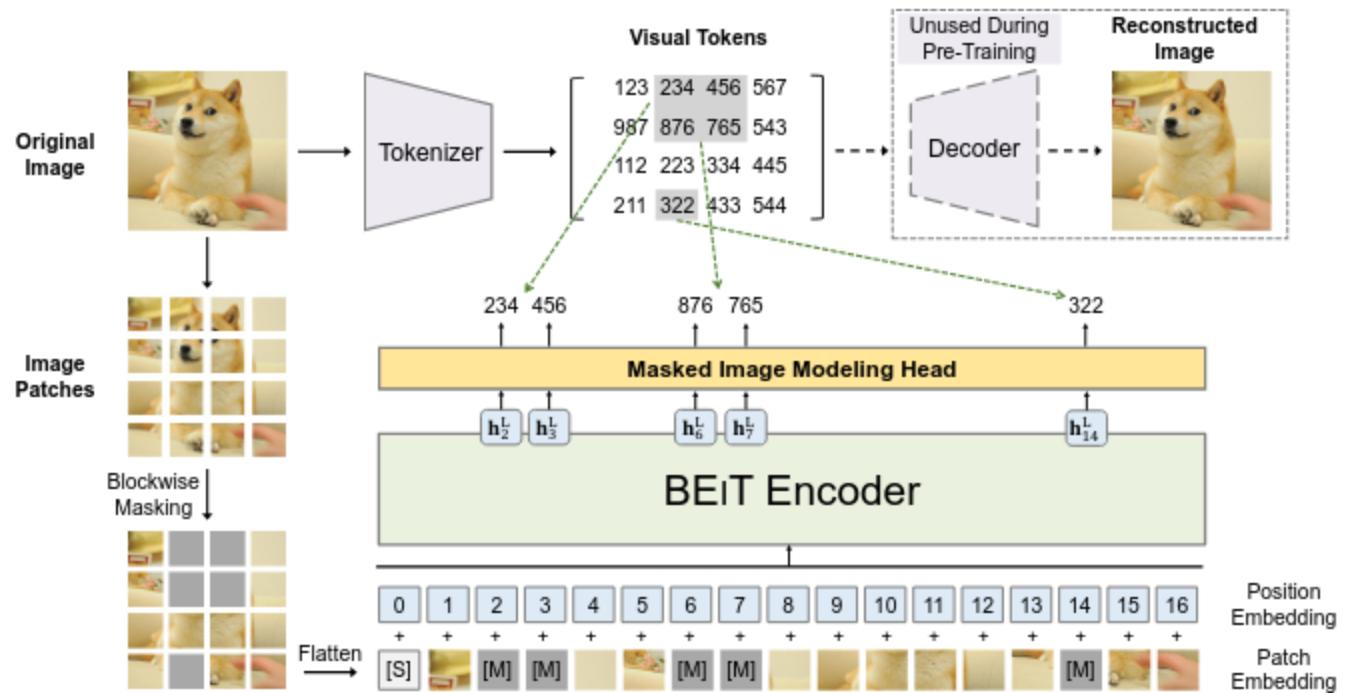
VISION TRANSFORMERS

Vision Transformers



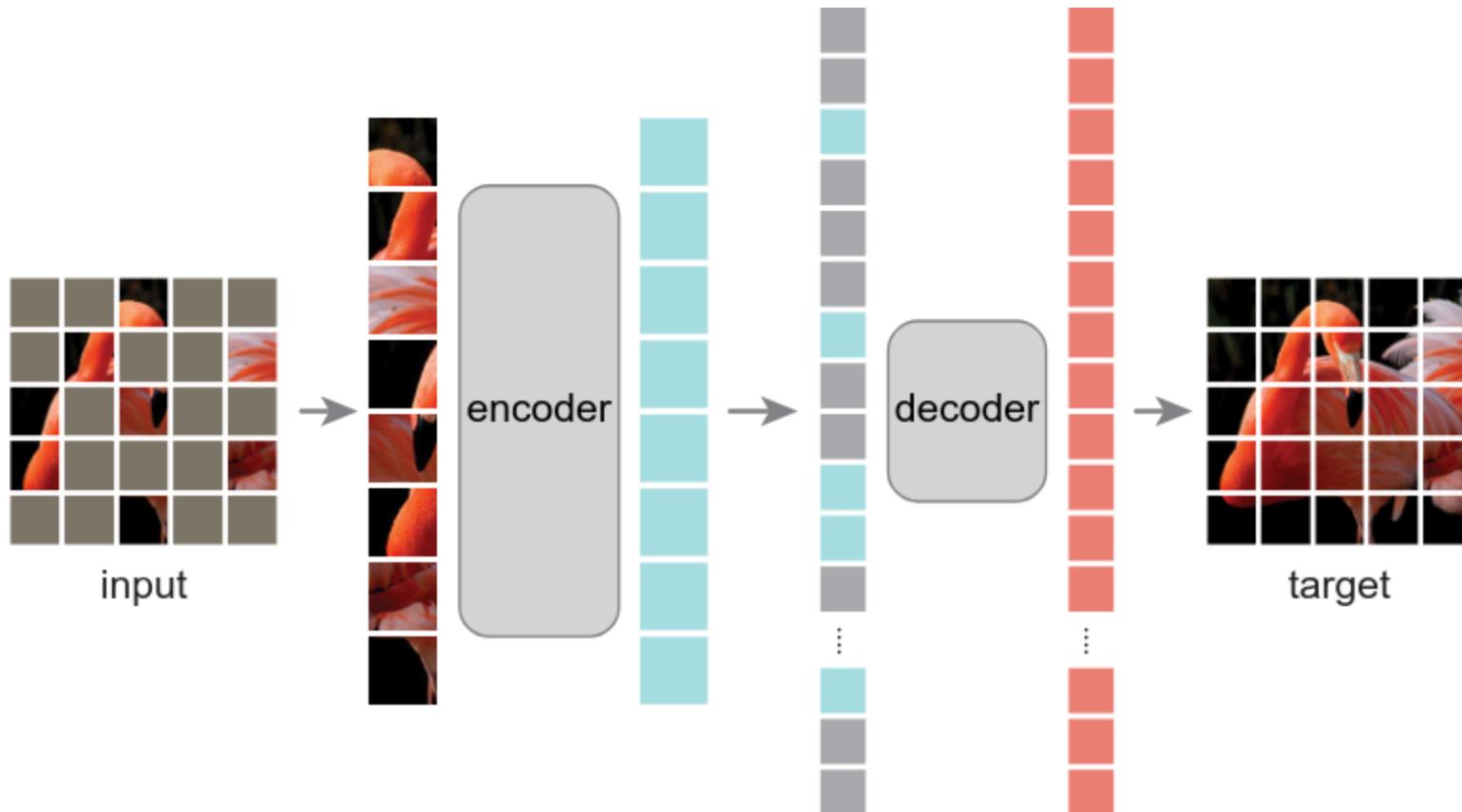
The NLP transformer architecture carries over directly to images: just change the **positional embedding**. Dosovitskiy et al. 2020, arXiv:2010.11929

BEiT - BERT Pre-Training of Image Transformers



BERT-like pretraining carries over directly. Bao et al., 2022

Masked Autoencoder (MAE)



MAE uses a simpler architecture and no tokenization. He et al. 2021

Masked Autoencoder - Reconstructions



He et al. 2021

Masked Autoencoder - Transfer Learning

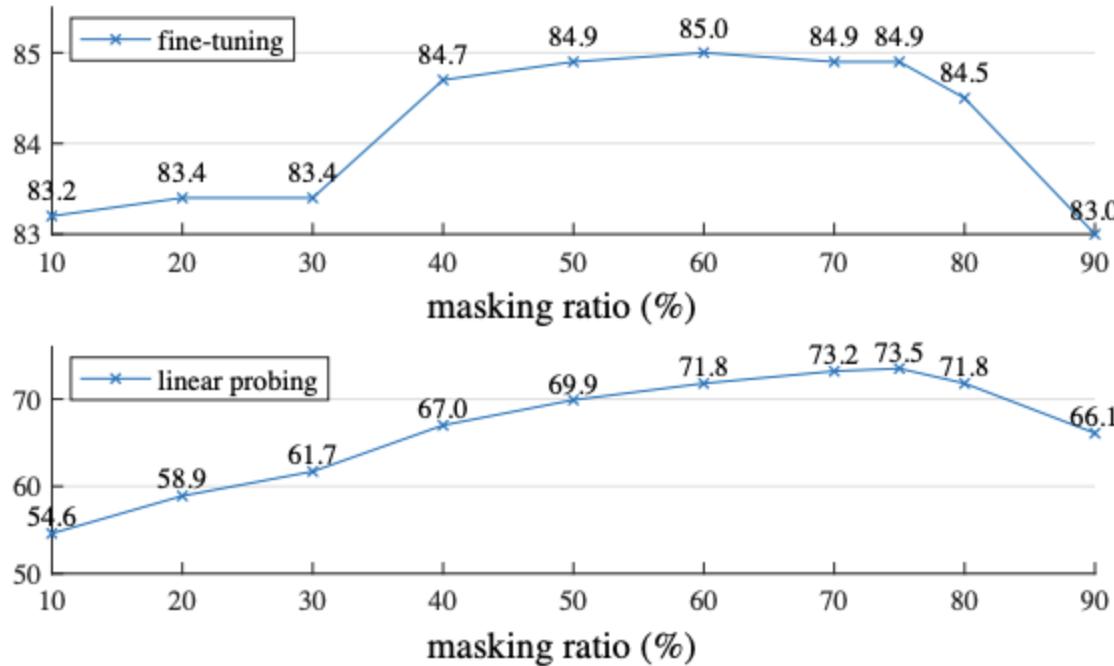


Figure 5. **Masking ratio.** A high masking ratio (75%) works well for both fine-tuning (top) and linear probing (bottom). The y-axes are ImageNet-1K validation accuracy (%) in all plots in this paper.

Unsupervised Training with Transformers

Transformer architectures make BERT-like masking useful unsupervised pre-training for image-related tasks.

OTHER APPROACHES

SimCLR - Contrastive Learning

Basic idea:

- generate two differently augmented versions of the same image
- train a representation that is as similar as possible for the same image, different for different images

Details:

- carefully choose augmentations, watch out for trivial solutions (e.g. color)
- separate representation from scoring
- compute "softmax over cosine similarity" over very large batches
- implemented with ResNet

SimCLR - Contrastive Learning

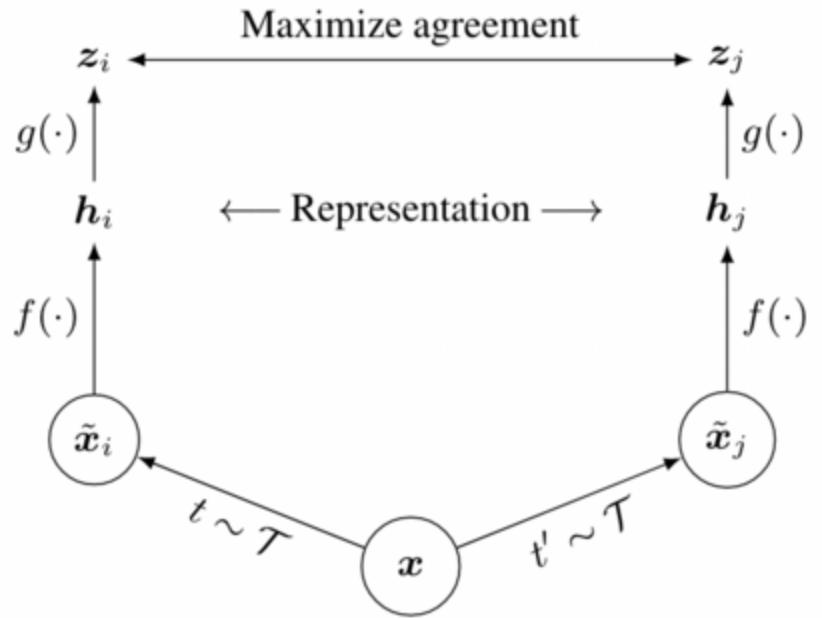


Figure 2. A simple framework for contrastive learning of visual representations. Two separate data augmentation operators are sampled from the same family of augmentations ($t \sim \mathcal{T}$ and $t' \sim \mathcal{T}$) and applied to each data example to obtain two correlated views. A base encoder network $f(\cdot)$ and a projection head $g(\cdot)$ are trained to maximize agreement using a contrastive loss. After training is completed, we throw away the projection head $g(\cdot)$ and use encoder $f(\cdot)$ and representation \mathbf{h} for downstream tasks.

SimCLR - Augmentations



(a) Original



(b) Crop and resize



(c) Crop, resize (and flip)



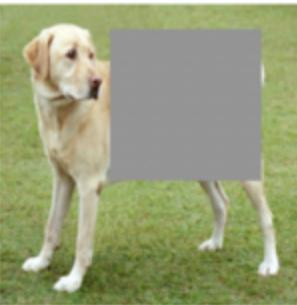
(d) Color distort. (drop)



(e) Color distort. (jitter)



(f) Rotate $\{90^\circ, 180^\circ, 270^\circ\}$



(g) Cutout



(h) Gaussian noise



(i) Gaussian blur



(j) Sobel filtering

Figure 4. Illustrations of the studied data augmentation operators. Each augmentation can transform data stochastically with some internal parameters (e.g. rotation degree, noise level). Note that we *only* test these operators in ablation, the *augmentation policy used to train our models* only includes *random crop (with flip and resize)*, *color distortion*, and *Gaussian blur*. (Original image cc-by: Von.grzanka)

SimCLR - Transfer Learning Performance

	Food	CIFAR10	CIFAR100	Birdsnap	SUN397	Cars	Aircraft	VOC2007	DTD	Pets	Caltech-101	Flowers
<i>Linear evaluation:</i>												
SimCLR (ours)	76.9	95.3	80.2	48.4	65.9	60.0	61.2	84.2	78.9	89.2	93.9	95.0
Supervised	75.2	95.7	81.2	56.4	64.9	68.8	63.8	83.8	78.7	92.3	94.1	94.2
<i>Fine-tuned:</i>												
SimCLR (ours)	89.4	98.6	89.0	78.2	68.1	92.1	87.0	86.6	77.8	92.1	94.1	97.6
Supervised	88.7	98.3	88.7	77.8	67.0	91.4	88.0	86.5	78.8	93.2	94.2	98.0
Random init	88.3	96.0	81.9	77.0	53.7	91.3	84.8	69.4	64.1	82.7	72.5	92.5

Table 8. Comparison of transfer learning performance of our self-supervised approach with supervised baselines across 12 natural image classification datasets, for ResNet-50 ($4\times$) models pretrained on ImageNet. Results not significantly worse than the best ($p > 0.05$, permutation test) are shown in bold. See Appendix B.8 for experimental details and results with standard ResNet-50.

DINO

- self-DIstillation with NO labels
- discovers labels / class structure by itself
- unsupervised representation learning for images
- impressive semantic segmentation results
- attention map = segmentation map
- vision transformer or ResNet 50 based

DINO Architecture

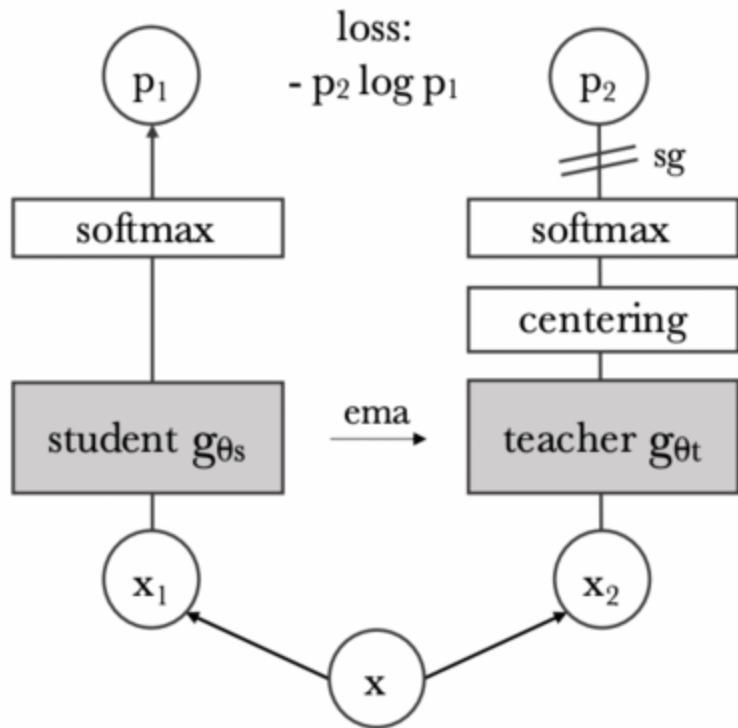


Figure 2: **Self-distillation with no labels.** We illustrate DINO in the case of one single pair of views (x_1, x_2) for simplicity. The

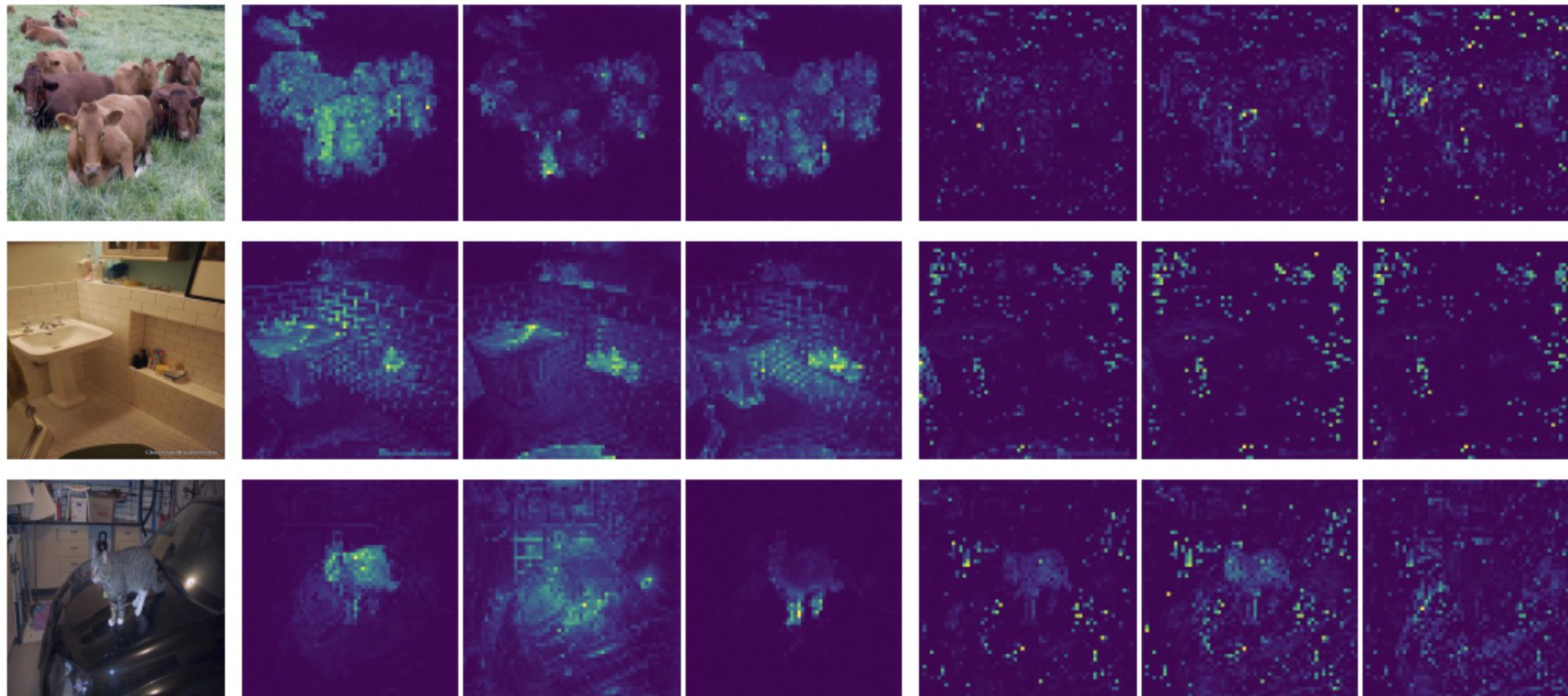
DINO Results

Method	Arch.	Param.	im/s	Linear	k -NN
Supervised	RN50	23	1237	79.3	79.3
SCLR [12]	RN50	23	1237	69.1	60.7
MoCov2 [15]	RN50	23	1237	71.1	61.9
InfoMin [67]	RN50	23	1237	73.0	65.3
BarlowT [81]	RN50	23	1237	73.2	66.0
OBoW [27]	RN50	23	1237	73.8	61.9
BYOL [30]	RN50	23	1237	74.4	64.8
DCv2 [10]	RN50	23	1237	75.2	67.1
SwAV [10]	RN50	23	1237	75.3	65.7
DINO	RN50	23	1237	75.3	67.5
Supervised	ViT-S	21	1007	79.8	79.8
BYOL* [30]	ViT-S	21	1007	71.4	66.6
MoCov2* [15]	ViT-S	21	1007	72.7	64.4
SwAV* [10]	ViT-S	21	1007	73.5	66.3
DINO	ViT-S	21	1007	77.0	74.5

DINO Segmentation

DINO

Supervised



Summary: SimCLR and DINO

High Level:

- SimCLR is a kind of *representation* or *metric* learning
- DINO is a kind of *clustering*
- implementations are complex and with lots of hyperparameters

Conclusions:

- tasks like these may be most useful as additional tasks combined with masking
(just like BERT)

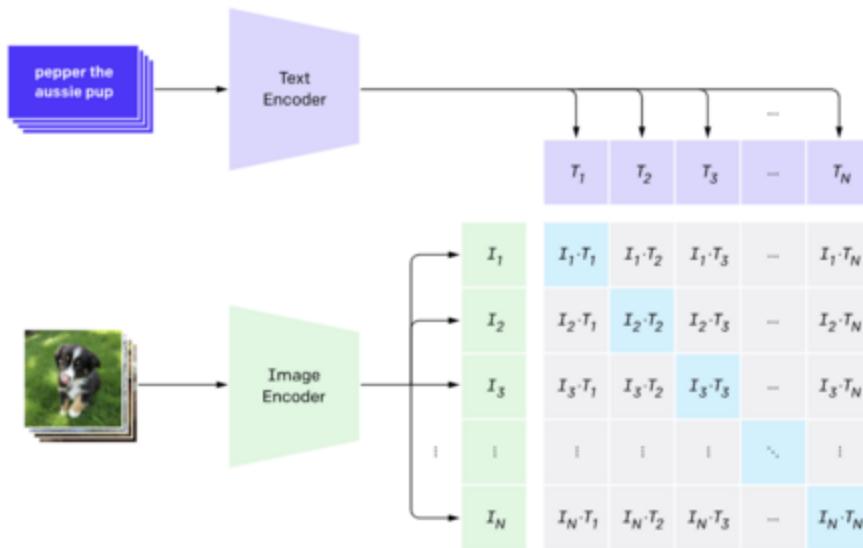
COMBINING TEXT AND IMAGE MODELS

Combining Image and Text Models

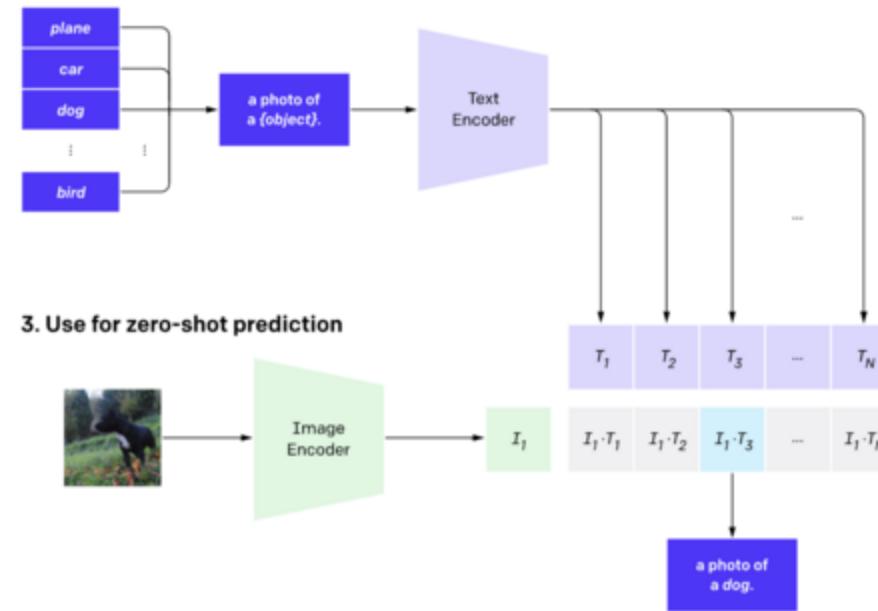
- GPT-3, ExT5, etc. show how natural language models can be used for zero shot learning
- CLIP
 - use natural language supervision for image recognition (weak supervision)
 - vision: transformer or ResNet, language: transformer
 - permit "prompt engineering" to allow different kinds of NLP tasks
 - uses contrastive pretraining (rather than, say, captioning)

CLIP Architecture

1. Contrastive pre-training



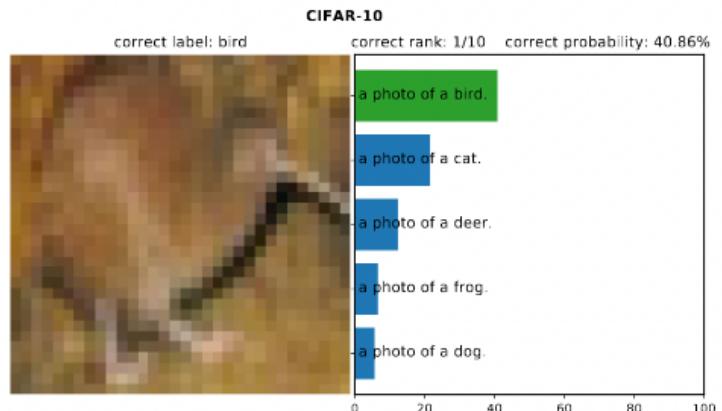
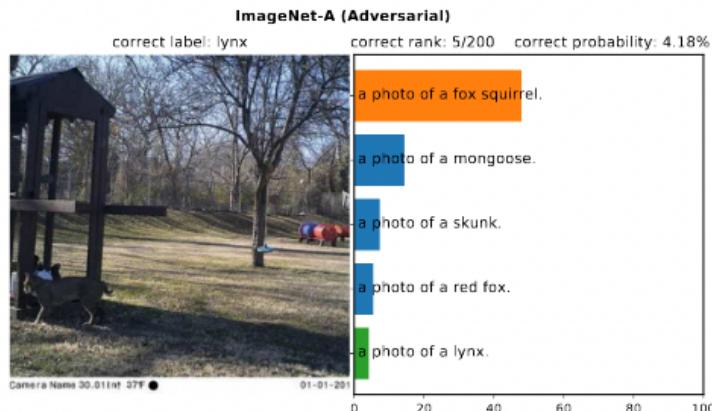
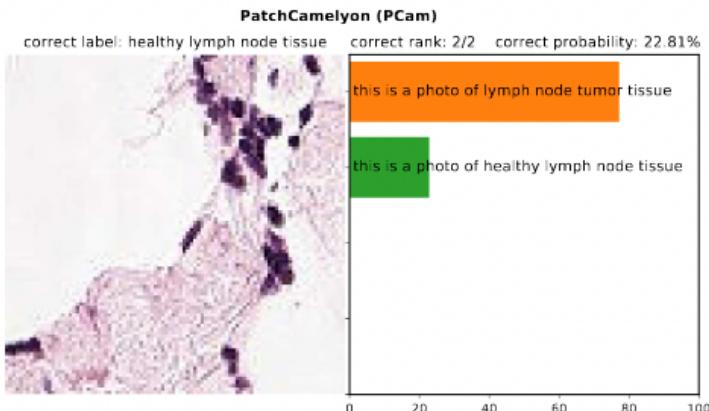
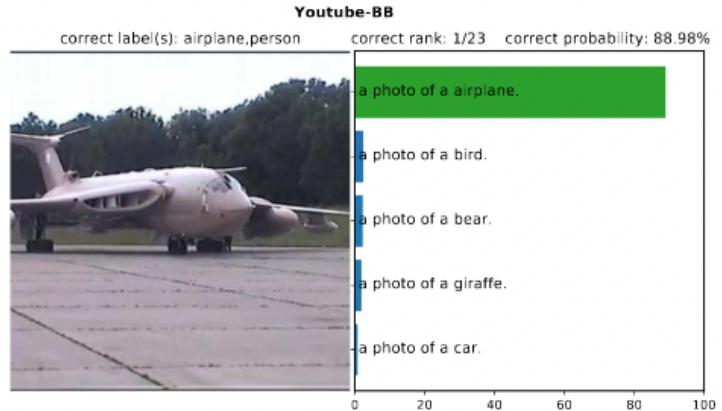
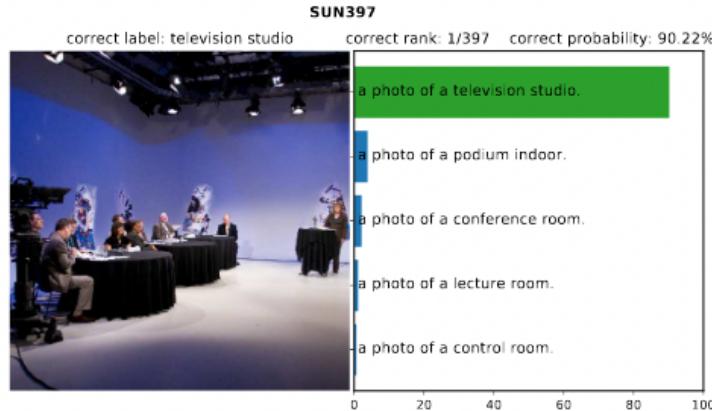
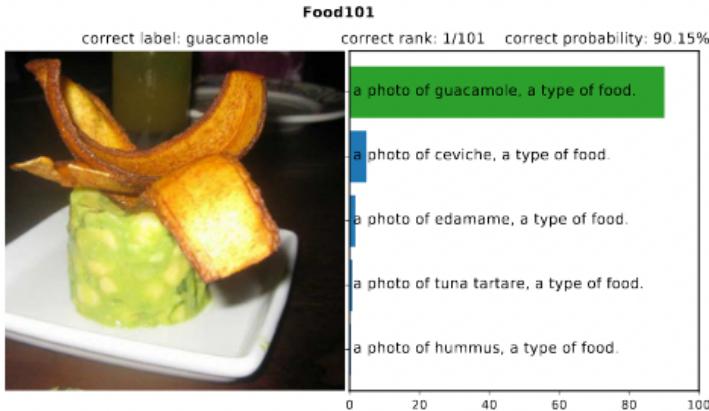
2. Create dataset classifier from label text



3. Use for zero-shot prediction



CLIP Results



CLIP Results

Zero-shot ImageNet accuracy

40%

20%

0%

2M

33M

67M

134M

268M

400M

Images processed

4x

3x efficiency

Bag of Words
Contrastive (CLIP)

Bag of Words
Prediction

Transformer
Language Model

Discussion

Current and future directions:

- combining text and image
- integrating unsupervised pretrained vision and language models
- integrating video and audio and identifying good self-supervised tasks for these (e.g., VideoMAE, Audio-MAE)
- cross-modal combinations likely also reduce the amount of training data required within each modality