

# Self-Supervision

Saurabh Gupta

# Is semantic supervision necessary to learn good representations?

- Manual labeling doesn't scale, suffers from biases
- Plenty of unlabeled visual data already, and growing really fast
- *And subject of the Gelato Bet:*
- *If, by the first day of autumn (Sept 23) of 2015, a method will exist that can match or beat the performance of R-CNN on Pascal VOC detection, without the use of any extra, human annotations (e.g. ImageNet) as pre-training, Mr. Malik promises to buy Mr. Efros one (1) gelato (2 scoops: one chocolate, one vanilla).*



# The Transformer: Transfer Learning

“ImageNet Moment for Natural Language Processing”

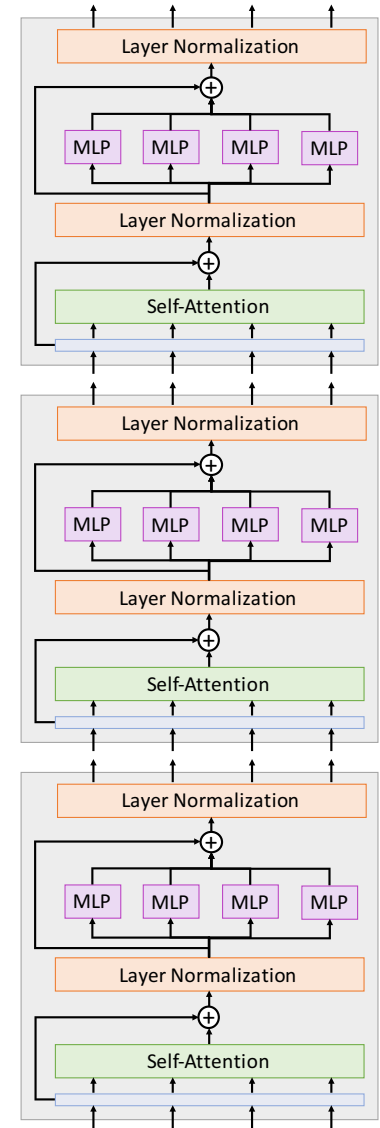
## Pretraining:

Download a lot of text from the internet

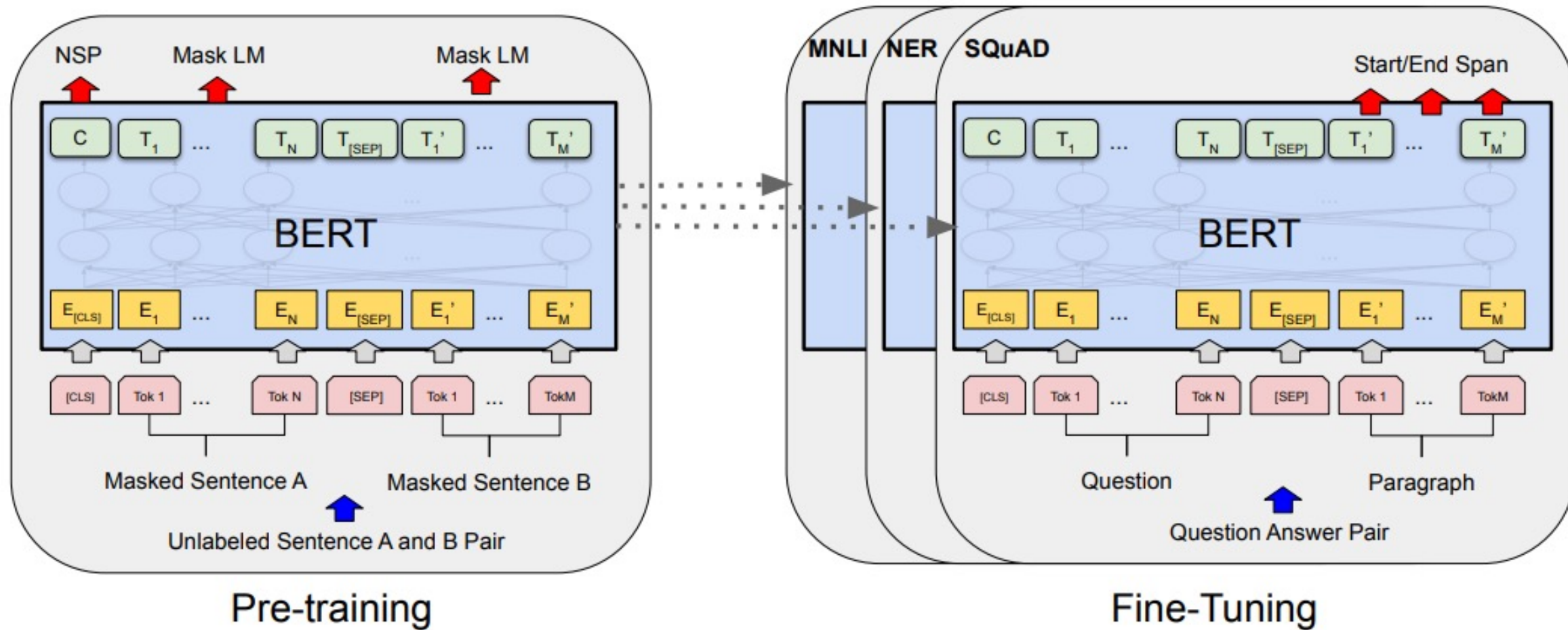
Train a giant Transformer model for language modeling

## Finetuning:

Fine-tune the Transformer on your own NLP task



# The Transformer: Transfer Learning



# The Transformer: Transfer Learning

System	MNLI-(m/mm) 392k	QQP 363k	QNLI 108k	SST-2 67k	CoLA 8.5k	STS-B 5.7k	MRPC 3.5k	RTE 2.5k	Average
Pre-OpenAI SOTA	80.6/80.1	66.1	82.3	93.2	35.0	81.0	86.0	61.7	74.0
BiLSTM+ELMo+Attn	76.4/76.1	64.8	79.8	90.4	36.0	73.3	84.9	56.8	71.0
OpenAI GPT	82.1/81.4	70.3	87.4	91.3	45.4	80.0	82.3	56.0	75.1
BERT <sub>BASE</sub>	84.6/83.4	71.2	90.5	93.5	52.1	85.8	88.9	66.4	79.6
BERT <sub>LARGE</sub>	<b>86.7/85.9</b>	<b>72.1</b>	<b>92.7</b>	<b>94.9</b>	<b>60.5</b>	<b>86.5</b>	<b>89.3</b>	<b>70.1</b>	<b>82.1</b>

Table 1: GLUE Test results, scored by the evaluation server (<https://gluebenchmark.com/leaderboard>). The number below each task denotes the number of training examples. The “Average” column is slightly different than the official GLUE score, since we exclude the problematic WNLI set.<sup>8</sup> BERT and OpenAI GPT are single-model, single task. F1 scores are reported for QQP and MRPC, Spearman correlations are reported for STS-B, and accuracy scores are reported for the other tasks. We exclude entries that use BERT as one of their components.

# Pre-train representations on a pre-text task

E.g. Colorization



After pre-training, use representation for down-stream tasks.

Many other possibilities,

- Spatial relationship between pair of patches

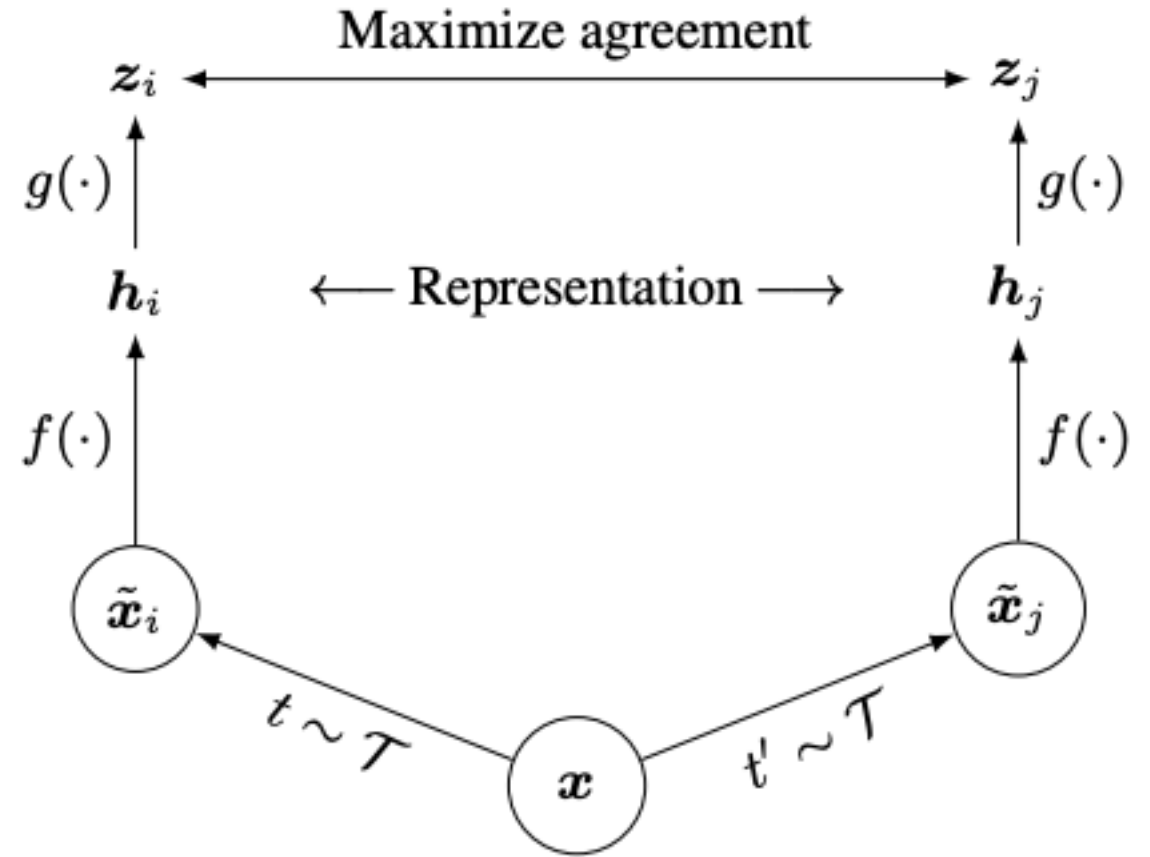


- Predict sound / frame ordering in a video
- Encourage two augmentations of same image to be closer to each other than to another image
- Predict hidden image patches from context

# Contrastive Learning

- Encourage two augmentations of an image to be close.
- Using a contrastive loss:

$$l_{i,j} = -\log \frac{\exp(\text{sim}(\mathbf{z}_i, \mathbf{z}_j)/\tau)}{\sum_{k=1}^{2N} \mathbb{1}_{[k \neq i]} \exp(\text{sim}(\mathbf{z}_i, \mathbf{z}_k)/\tau)}$$



[A Simple Framework for Contrastive Learning of Visual Representations](#), Chen et al. ICML 2020

See also: [Momentum Contrast for Unsupervised Visual Representation Learning](#), He et al. CVPR 2020

# 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



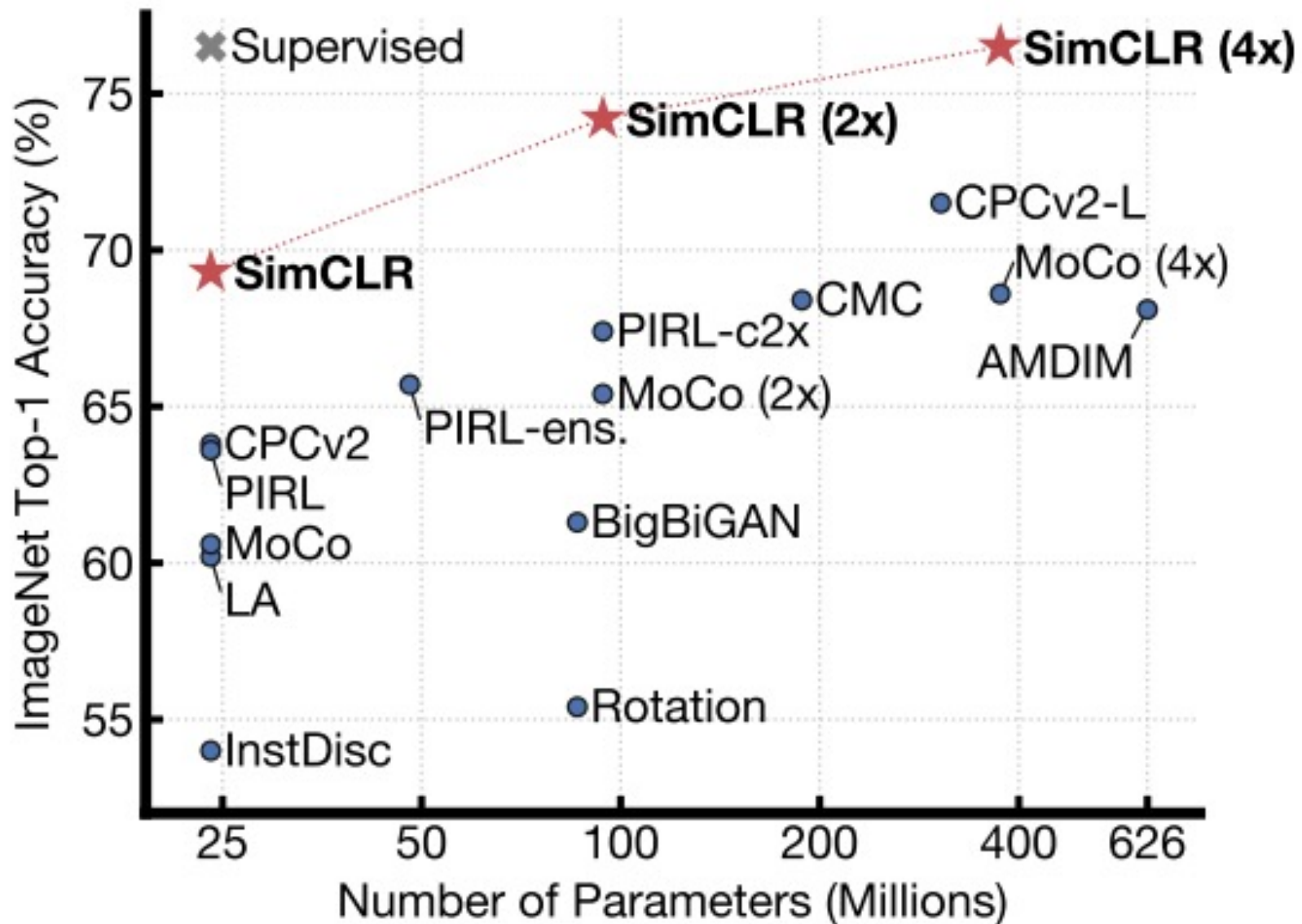
(i) Gaussian blur



(j) Sobel filtering

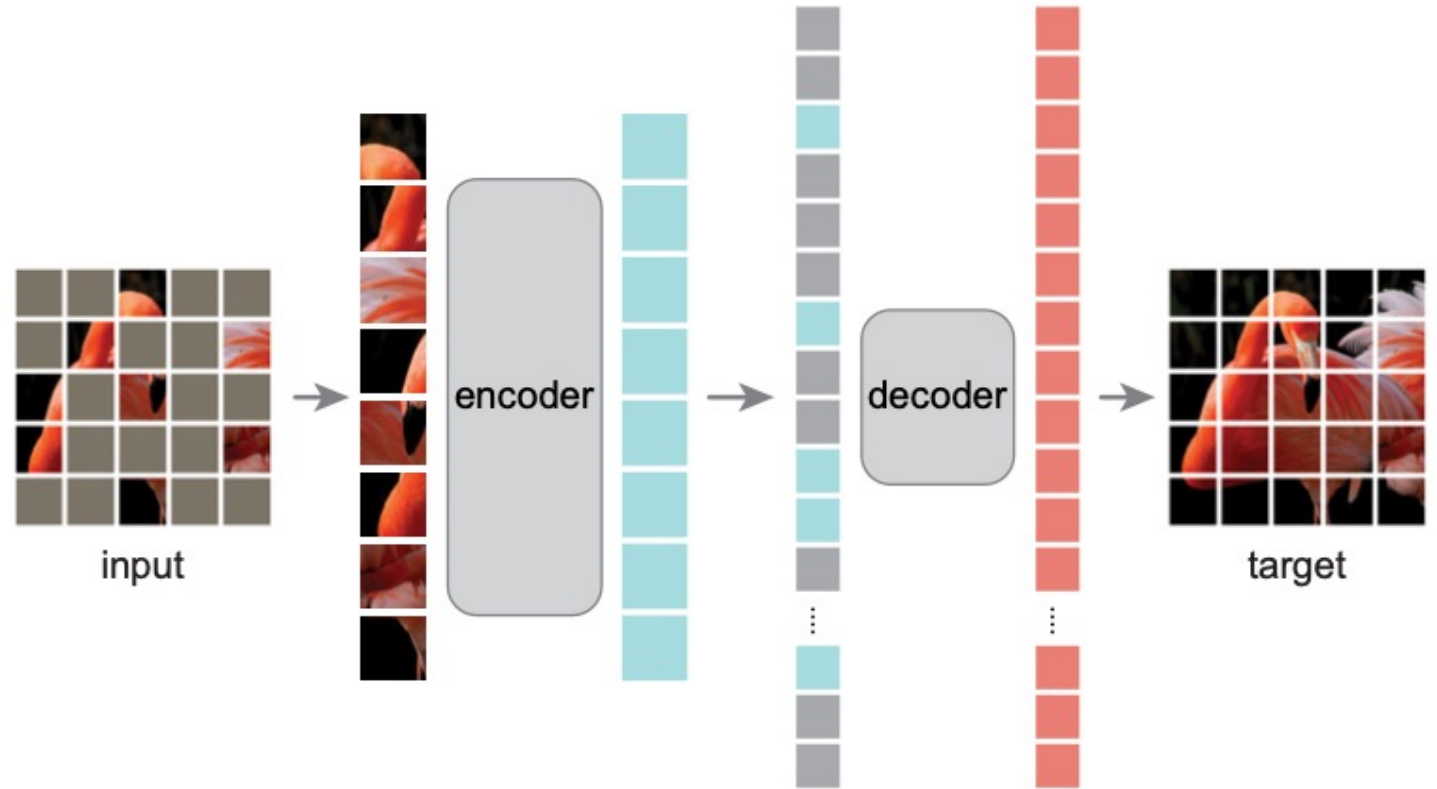


# Results



# Masked Auto-Encoders

- Mask out image patches, predict masked patches from visible patches.
- Pre-train encoder & decoder.
- Use encoder as an image representation.



# Better than semantic supervision on ImageNet 1K!

method	pre-train data	AP <sup>box</sup>		AP <sup>mask</sup>	
		ViT-B	ViT-L	ViT-B	ViT-L
supervised	IN1K w/ labels	47.9	49.3	42.9	43.9
MoCo v3	IN1K	47.9	49.3	42.7	44.0
BEiT	IN1K+DALLE	49.8	<b>53.3</b>	44.4	47.1
MAE	IN1K	<b>50.3</b>	<b>53.3</b>	<b>44.9</b>	<b>47.2</b>

Table 4. **COCO object detection and segmentation** using a ViT Mask R-CNN baseline. All entries are based on our implementation. Self-supervised entries use IN1K data *without* labels. Mask AP follows a similar trend as box AP.

method	pre-train data	ViT-B	ViT-L
supervised	IN1K w/ labels	47.4	49.9
MoCo v3	IN1K	47.3	49.1
BEiT	IN1K+DALLE	47.1	53.3
MAE	IN1K	<b>48.1</b>	<b>53.6</b>

Table 5. **ADE20K semantic segmentation** (mIoU) using UperNet. BEiT results are reproduced using the official code. Other entries are based on our implementation. Self-supervised entries use IN1K data *without* labels.

dataset	ViT-B	ViT-L	ViT-H	ViT-H <sub>448</sub>	prev best
iNat 2017	70.5	75.7	79.3	<b>83.4</b>	75.4 [55]
iNat 2018	75.4	80.1	83.0	<b>86.8</b>	81.2 [54]
iNat 2019	80.5	83.4	85.7	<b>88.3</b>	84.1 [54]
Places205	63.9	65.8	65.9	<b>66.8</b>	66.0 [19] <sup>†</sup>
Places365	57.9	59.4	59.8	<b>60.3</b>	58.0 [40] <sup>‡</sup>

Table 6. **Transfer learning accuracy on classification datasets**, using MAE pre-trained on IN1K and then fine-tuned. We provide system-level comparisons with the previous best results.

<sup>†</sup>: pre-trained on 1 billion images. <sup>‡</sup>: pre-trained on 3.5 billion images.

# Improves performance on ImageNet itself

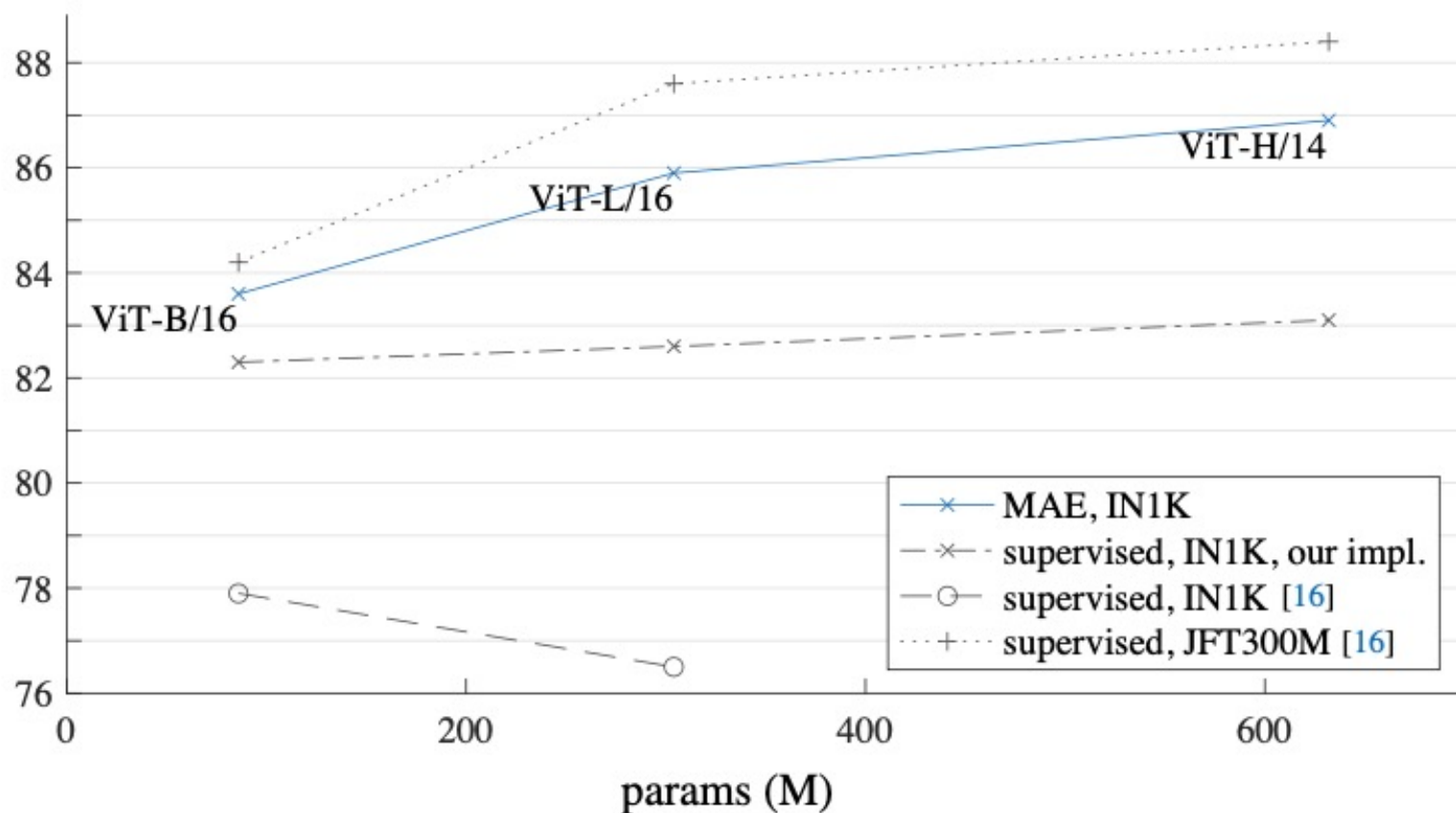
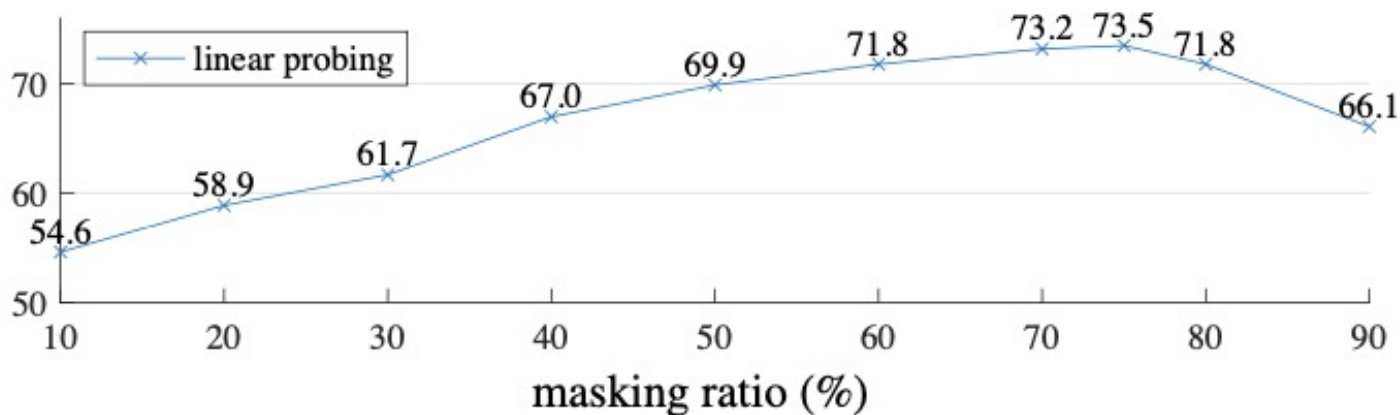
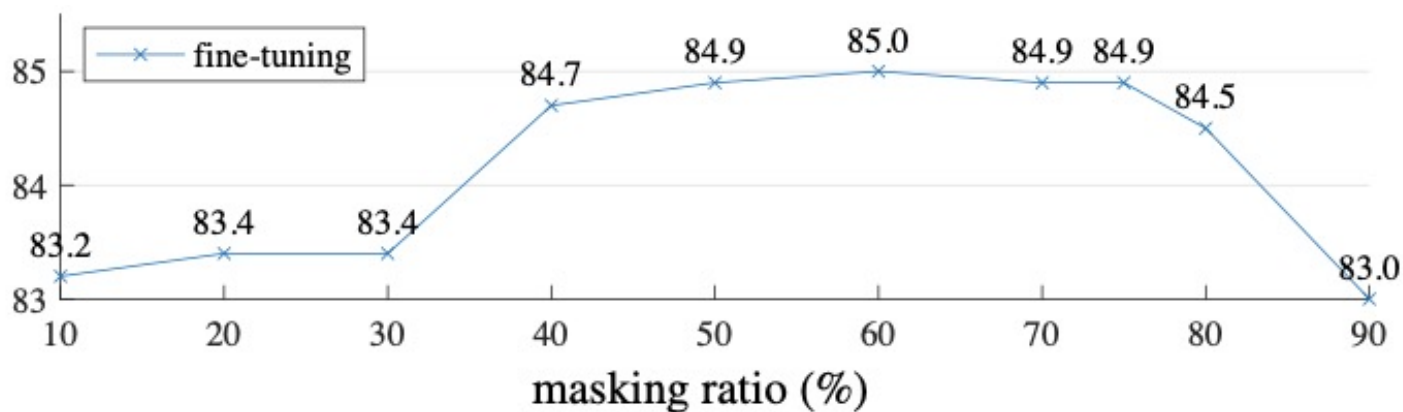


Figure 8. **MAE pre-training vs. supervised pre-training**, evaluated by fine-tuning in ImageNet-1K (224 size). We compare with the original ViT results [16] trained in IN1K or JFT300M.

# Better than past self-supervision approaches

method	pre-train data	ViT-B	ViT-L	ViT-H	ViT-H <sub>448</sub>
scratch, our impl.	-	82.3	82.6	83.1	-
DINO [5]	IN1K	82.8	-	-	-
MoCo v3 [9]	IN1K	83.2	84.1	-	-
BEiT [2]	IN1K+DALLE	83.2	85.2	-	-
MAE	IN1K	<u>83.6</u>	<u>85.9</u>	<u>86.9</u>	<b>87.8</b>

# Ablations



Need high masking ratio for good learning.  
NLP models use 15-20% masking ratio.

case	ft	lin	FLOPs
encoder w/ [M]	84.2	59.6	3.3×
encoder w/o [M]	<b>84.9</b>	<b>73.5</b>	<b>1×</b>

Faster and better to not input masked out patches to encoder

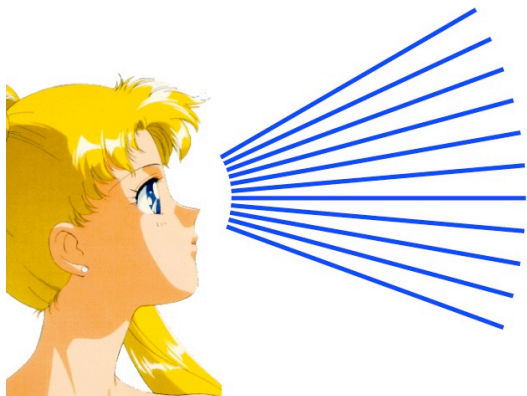
case	ft	lin
pixel (w/o norm)	84.9	73.5
pixel (w/ norm)	<b>85.4</b>	<b>73.9</b>
PCA	84.6	72.3
dVAE token	85.3	71.6

Normalized pixels are a better target than discrete tokens / PCA coefficients

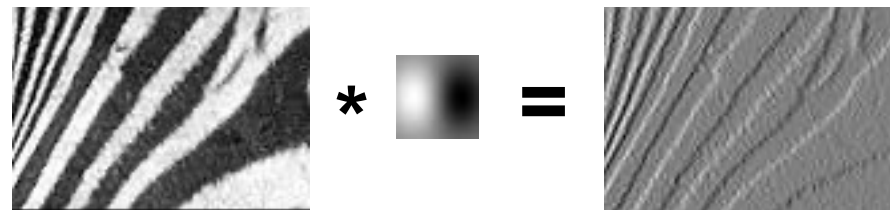
# I. Early vision

---

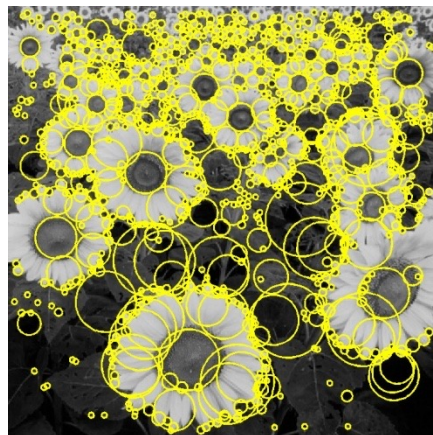
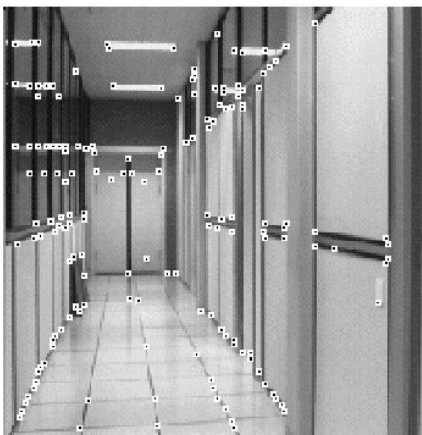
## Basic image formation and processing



Cameras and sensors  
Light and color



Linear filtering  
Edge detection



Feature extraction

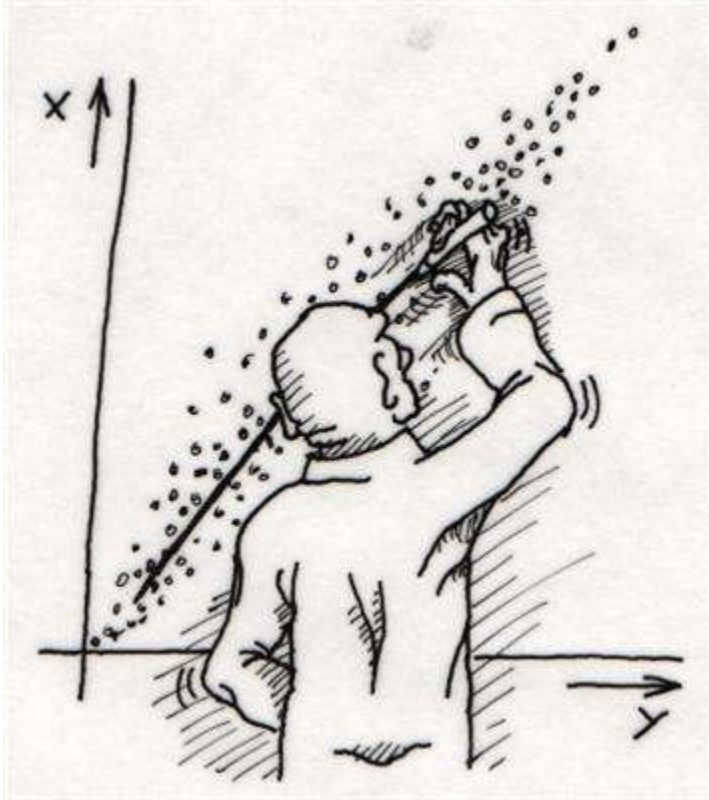


Optical flow

## II. “Mid-level vision”

---

### Fitting and grouping



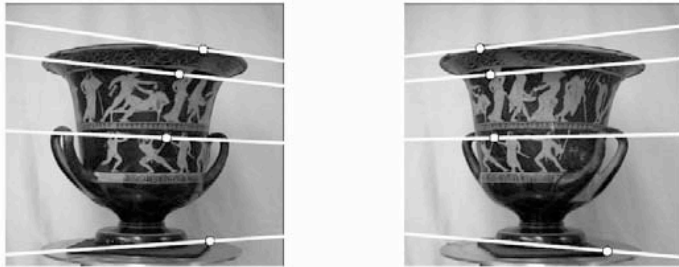
Fitting: Least squares  
Voting methods



Alignment



# III. Multi-view geometry



Epipolar geometry

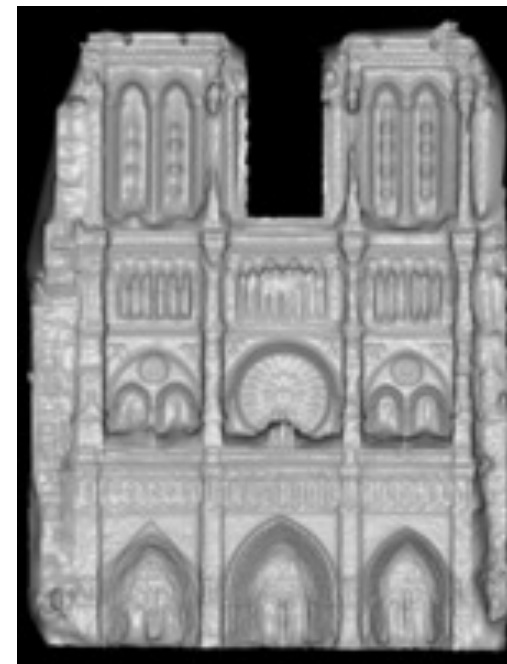


Two-view stereo



Драконъ, видимый подъ различными углами зрѣнія  
По гравюру на мѣди изъ „Oculus artificialis teleiopicus“ Цана. 1702 года.

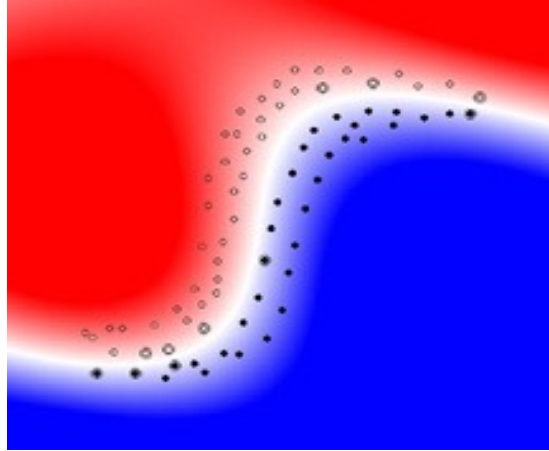
Structure from motion



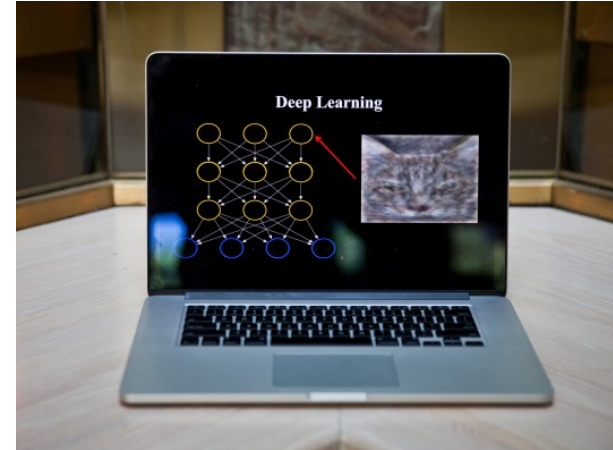
Multi-view stereo

# IV. Recognition

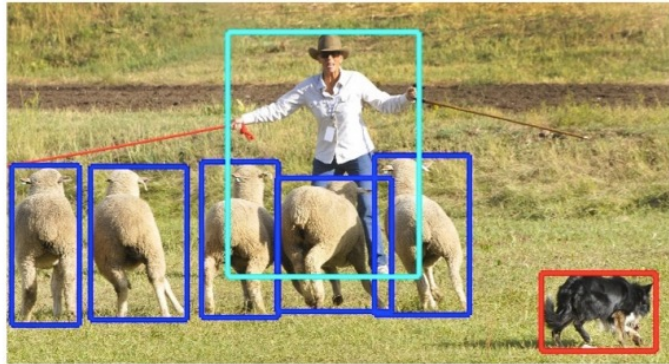
---



Basic classification



Deep learning



Object detection



Segmentation

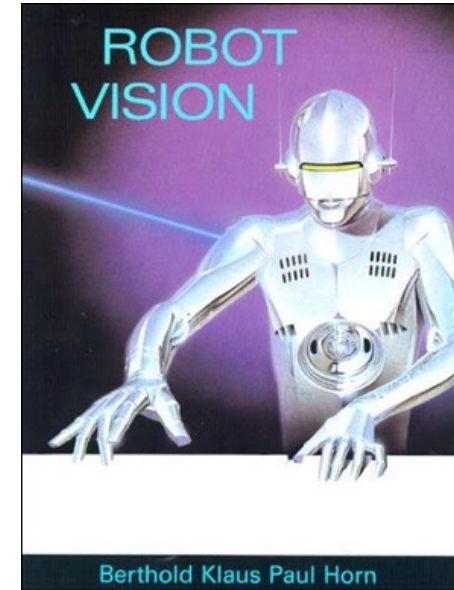
# V. Additional Topics (time permitting)



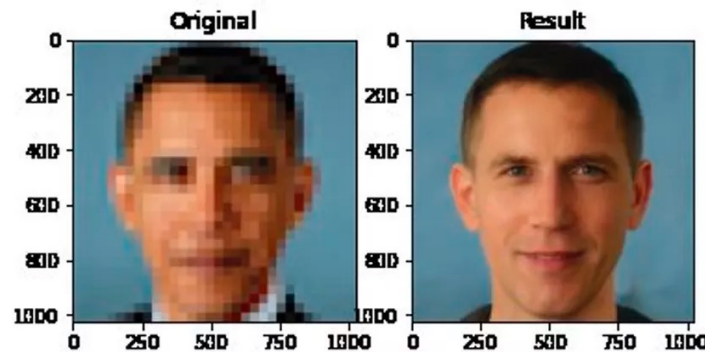
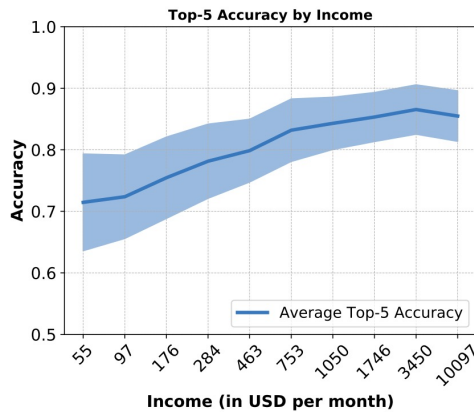
Video



Generation



Vision and Robotics



Bias and Ethical Considerations



# What's next?

- [Machine Learning \(CS 446\)](#)
- [Applied Machine Learning \(CS 441\)](#)
- [Deep Learning for Computer Vision \(CS 444\)](#)
- **Advanced Classes:**
  - [Robot Perception](#) (Shenlong)
  - [3D Vision](#) (Derek)
  - [Robot Learning](#) (Saurabh, Yunzhu)
  - [Autonomous Vehicles](#) (DAF)
  - [Learning to Learn](#) (Yuxiong)
  - Efficient & Predictive Vision (Lynna)
  - [Meta-Vision](#) (Lana)
  - [Deep Generative and Dynamical Models](#) (Arindam)
  - [Generative AI Models](#) (Lav)