

# Unsupervised Pre-Training: Contrastive Learning

CS 330

# Course Reminders

Project proposal due Wednesday.  
(graded lightly, for your benefit)

Homework 2 due next Monday 10/24.

Following up on some high-res feedback:

- I will work on making whiteboard writing larger.
- Moving one TA office hours (Garrett) from in-person-> over zoom.
- Will clarify project expectations on Ed.

# So Far

Few-shot learning via meta-learning

Problem: Given data from  $\mathcal{T}_1, \dots, \mathcal{T}_n$ , solve new task  $\mathcal{T}_{\text{test}}$  more quickly / proficiently / stably

Methods: black-box, optimization-based, non-parametric

What if you don't have a lot of tasks?

What if you *only* have one batch of unlabeled data?

# This Lecture

Unsupervised representation learning for few-shot learning

Part I: Contrastive learning

Part II (next time): Reconstruction-based methods

Relation to meta-learning.

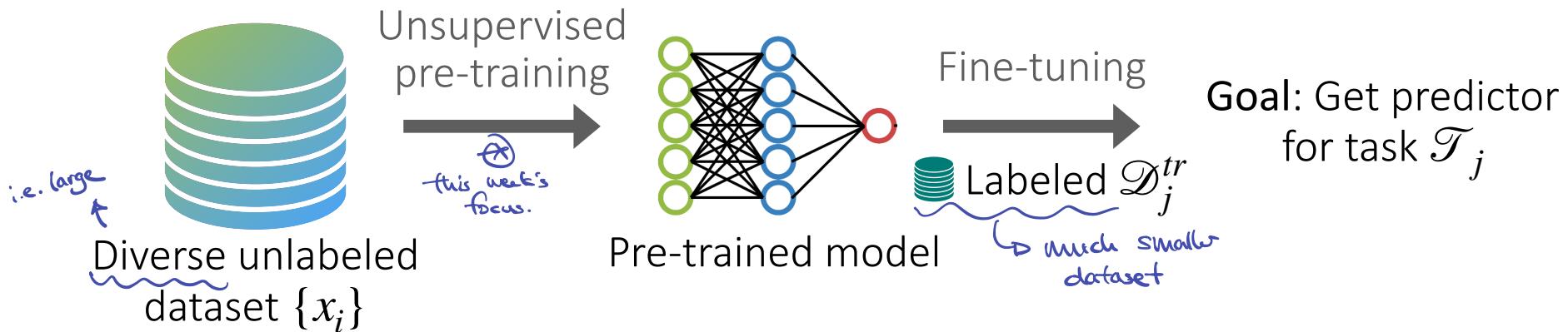
*pretty closely related!*

## Goals for the lecture:

- Understand **contrastive learning**: intuition, design choices, how to implement
- How contrastive learning relates to meta-learning

# Unsupervised Pre-Training Set-Up

Did this for point cloud data + SimCLR.



# Key Idea of Contrastive Learning

**Similar examples** should have **similar representations**

key idea of  
contrastive learning!

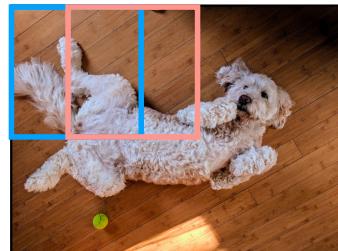
Examples with the same class label



(Requires labels, related to Siamese nets, ProtoNets)

↳ Covered  
last class

Nearby image patches



Augmented versions of the example



(flip & crop)

Diff types of  
similar images

Nearby video frames



van den Oord, Li, Vinyals. CPC. 2018

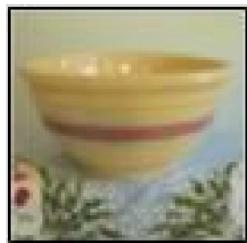
Chen, Kornblith, Norouzi, Hinton. SimCLR. ICML 2020

# Key Idea of Contrastive Learning

**Similar examples should have similar representations**



Similar representations



Similar representations

~~Important~~ Important

Question: Why not simply minimize difference between representations?

$$\min_{\theta} \sum_{(x_i, x_j)} \|f_{\theta}(x_i) - f_{\theta}(x_j)\|^2$$

Need to both compare & contrast!

→ Converges to a degenerate solution where all images get mapped to a single constant vector, i.e.  $\mathcal{L} \rightarrow 0$

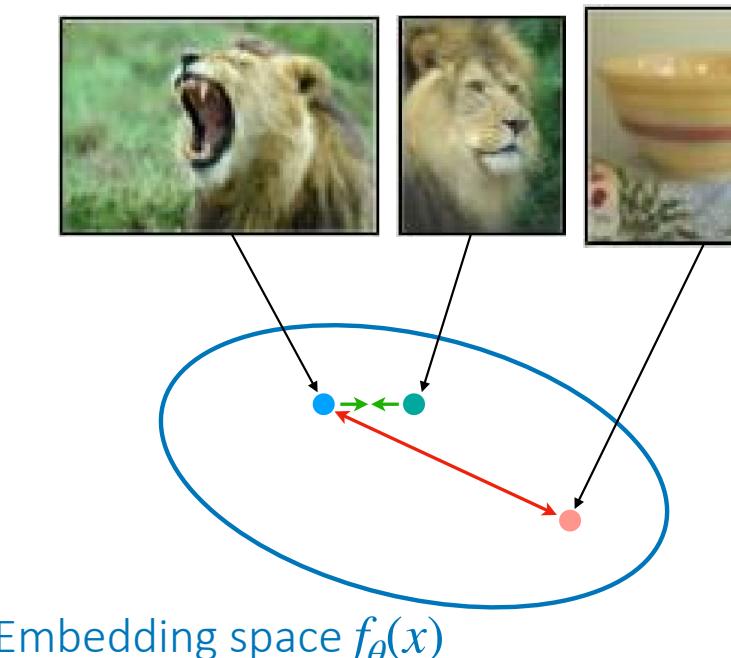
⇒ Don't get a good representation... just collapses.

⇒ Hence the "Contrastive" learning

# Key Idea of Contrastive Learning

🚫 **Similar examples** should have **similar representations**

Need to both compare & contrast!



Bring together representations of similar examples.

🚫 Push apart representations of differing examples.

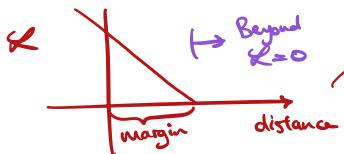
⭐ Key design choices:

1. Implementation of contrastive loss
2. Choosing what to compare/contrast

이거 뭔말이지?

# Contrastive Learning Implementation

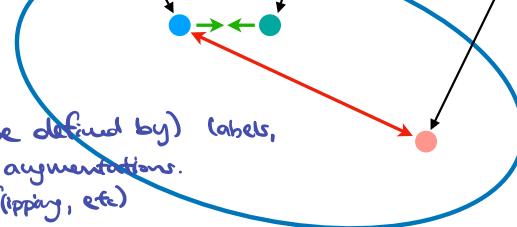
⊕ Hinge Loss (applied in the eqn.)



Similar examples should have similar representations



anchor  $x$  positive  $x^+$  negative  $x^-$



(+)ve ~ (-)ve  
can come from (be defined by) labels,  
or! can also use augmentations.  
(random crop, flipping, etc)

Embedding space  $f_\theta(x)$

not always exactly  
class labels

Need to both compare & contrast!

⊕ Simplest form

V1. Triplet loss:

$$\min_{\theta} \sum_{(x, x^+, x^-)} \max(0, \|f_\theta(x) - f_\theta(x^+)\|^2 - \|f_\theta(x) - f_\theta(x^-)\|^2 + \epsilon)$$

can we Euclidean  
or negative  
Cosine Similarity

⊕ this is unbounded  
⇒ will go to  $\leftarrow$ ve inf. unless  
you take this into account  
when designing embedding func.  
or the loss function.

Q: Is there a way to push some examples more than others?

A. In general CC pushes away the same for all types.  
Examples may end up in diff. distances ∵ some will be naturally be harder to push away.



Compare to Siamese networks:

Classify  $(x, x')$  as same class if  $\|f(x) - f(x')\|^2$  is small.

⊕ Key difference: learns a metric space, not just a classifier

↳ triplet loss

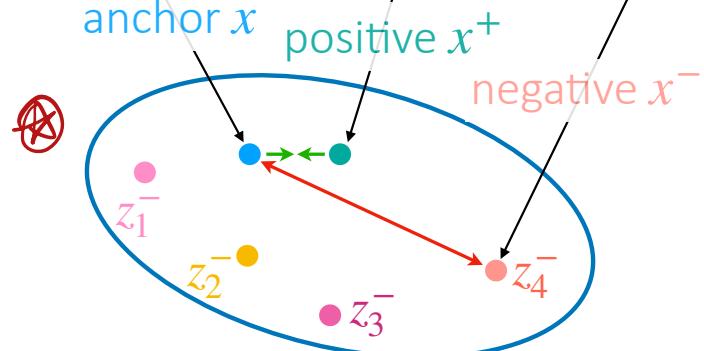
Challenge: need to find difficult negatives.

- If you sample a (+)ve that's already u. far away, then you'll just get  $\mathcal{L}=0$ , v the model will not learn.
- "Hard Negative" helps find (-)ves that allows the model to continually learn.

# Contrastive Learning Implementation

Similar examples should have similar representations

Need to both compare & contrast!



Embedding space  $f_\theta(x)$

Trying to classify whether an example is (+)ve or (-)ve, given an anchor (same w/ Vi-triplet loss)

V2. From binary to N-way classification:

↳ rather than binaries ↳ multiple negatives

$$\mathcal{L}_{\text{N-way}}(\theta) = - \sum_z \log \frac{\exp(-d(z, z^+))}{\sum_i \exp(-d(z, z_i^-))}$$

$z = f_\theta(x)$

↳ Used in SimCLR paper! - generalization of triplet loss to multiple negatives

↳ V. Similar to summing up the negative loss for triplet loss instead of  $\Sigma$ , it's doing  $\Sigma \log$ .

↳ Softmax ( $-d(\text{anchor}, \text{example})$ )  
↳ Could be both (+)ve & (-)ve

↳ Alternatively, can sum all examples in the denominator

# Contrastive Learning Implementation

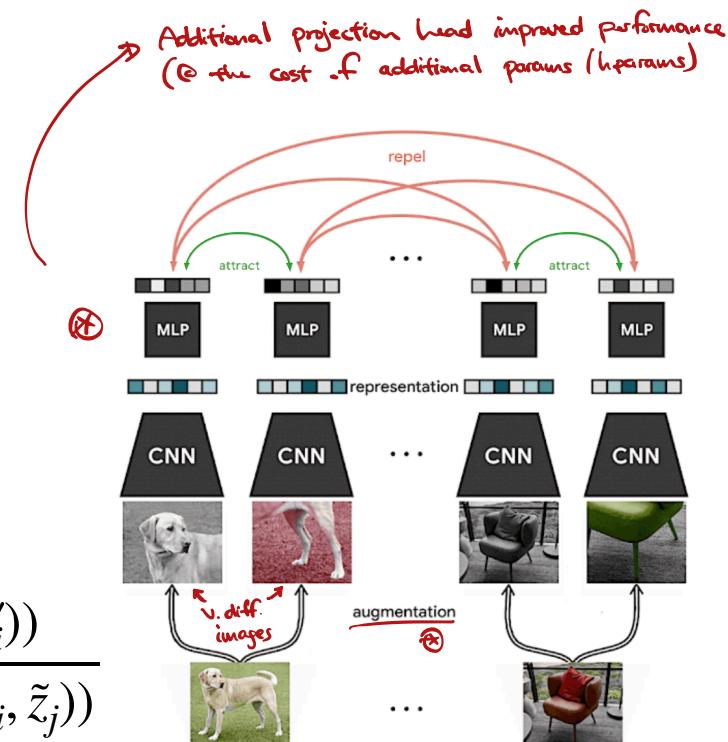
⚠ Generating good augmentations is key.

## SimCLR Algorithm

### Unsupervised Pre-Training

1. Sample minibatch of examples  $x_1, \dots, x_N$
2. Augment each example twice to get  $\tilde{x}_1, \dots, \tilde{x}_N, \tilde{x}'_1, \dots, \tilde{x}'_N$
3. Embed examples with  $f_\theta$  to get  $\tilde{z}_1, \dots, \tilde{z}_N, \tilde{z}'_1, \dots, \tilde{z}'_N$
4. Compute all pairwise distances  $d(z_i, z_j) = -\frac{z_i^T z_j}{\|z_i\| \|z_j\|}$
5. Update  $\theta$  w.r.t. loss  $\mathcal{L}_{\text{N-way}}(\theta) = -\sum_i \log \frac{\exp(-d(\tilde{z}_i, \tilde{z}'_i))}{\sum_{j \neq i} \exp(-d(\tilde{z}_i, \tilde{z}_j))}$

generates more ←'s  
can augment more or less.



**After Pre-Training:** train classifier on top of representation or fine-tune entire network.

⚠ Data imbalance may negatively affect the performance.

⇒ May be difficult to notice esp. for unlabeled datasets.

# Performance of Contrastive Learning

ImageNet Classification Results

Method	Architecture	Label fraction		
		1%	10%	Top 5
Supervised baseline	ResNet-50	48.4	80.4	
<i>Methods using other label-propagation:</i>				
Pseudo-label	ResNet-50	51.6	82.4	
VAT+Entropy Min.	ResNet-50	47.0	83.4	
UDA (w. RandAug)	ResNet-50	-	88.5	
FixMatch (w. RandAug)	ResNet-50	-	89.1	
S4L (Rot+VAT+En. M.)	ResNet-50 (4×)	-	91.2	
<i>Methods using representation learning only:</i>				
InstDisc	ResNet-50	39.2	77.4	
BigBiGAN	RevNet-50 (4×)	55.2	78.8	
PIRL	ResNet-50	57.2	83.8	
CPC v2	ResNet-161(*)	77.9	91.2	
SimCLR (ours)	ResNet-50	75.5	87.8	
SimCLR (ours)	ResNet-50 (2×)	83.0	91.2	
SimCLR (ours)	ResNet-50 (4×)	85.8	92.6	

Table 7. ImageNet accuracy of models trained with few labels.

Semi-Supervised learning

Used n% of the overall label data → train SimCLR → do clf

→ Almost in few-shot learning regime

1% labels: ~12.8 images/class

10% labels : ~128 images / class

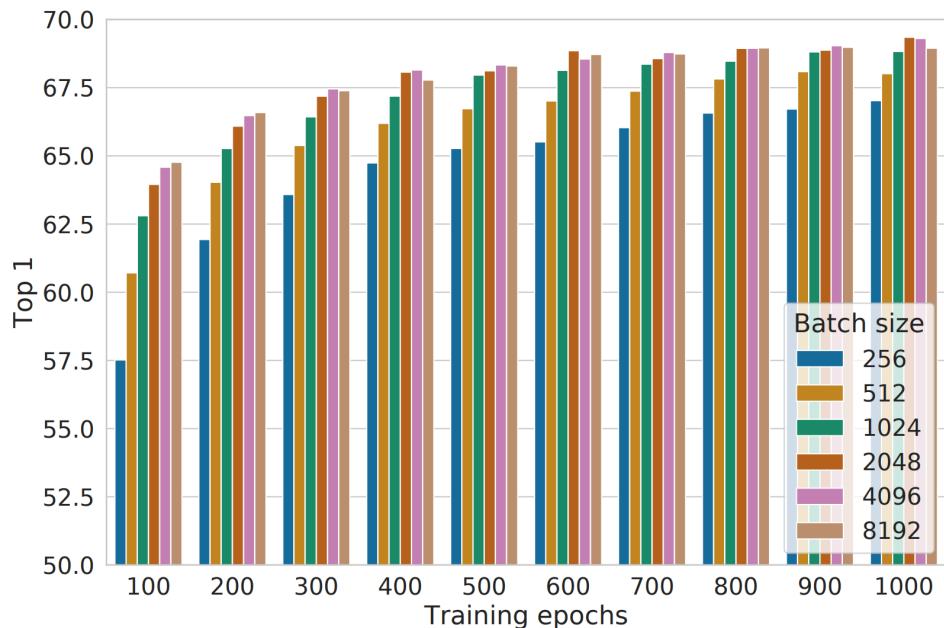
] but also uses the entire rest of the images (w/o labels)

- Substantial improvements over training from scratch
- Improvements over other methods, especially in 1% label setting

→ width of the hidden layer is 4x larger.

# Performance of Contrastive Learning

## Effect of Batch Size & Number of Training Epochs



- Important to train for longer ( $\sim 600+$  epochs)
- Requires large batch size

# Why does contrastive learning need a large batch size?

Interpretation of loss: classifying augmented example from rest of dataset

$$\mathcal{L}_{N\text{-way}}(\theta) = - \sum_i \log \frac{\exp(-d(\tilde{z}_i, \tilde{z}'_i))}{\sum_{j \neq i} \exp(-d(\tilde{z}_i, \tilde{z}_j))} \quad \leftarrow \text{summation over entire dataset}$$


Intuition: Closest  $z$  will dominate the denominator, can be missed when subsampling

Mathematically?

↳ applies just to  $(j \neq i)$  samples → samples from different class (for instance)

⊗ negatives w/ the smallest negative distance will dominate the sum

$$\begin{array}{c} \hookrightarrow d: \begin{array}{cc} \overset{j \neq i}{\text{similar}} & \overset{\text{dissimilar}}{10} \\ \overset{0.01}{e^{-0.01}} & e^{-10} \end{array} \\ \text{plays a much larger role} \\ (0.99) \quad (0.00004) \end{array}$$

Because these examples w/ small distances are dominating,  
if you're subsampling (have a small batch size),  
you can miss the examples that actually dominate that loss func.  
⇒ Need better representation of more dissimilar examples.

cf) In supervised learning, it's totally fine to have small batches  $\rightarrow -\sum_{x,y} \log P(y|x)$ . Smaller batches of  $x \sim y$  will still give the correct gradient.  $\rightarrow$  Not the same here.

$$\min_{\theta} -\log \frac{e^{-d(\bar{z}, \bar{z}^+)}}{\sum_n e^{-d(z, z_n)}}$$

↳ Subsampling: sort of bringing the summation outside of the log.  
 $\Rightarrow d(\bar{z}, \bar{z}^+) + \log \sum_n e^{-d(z, z_n)}$   
 when we subsample.  $\downarrow$  Summation is inside the log.  
 $\approx d(\bar{z}, \bar{z}^+) + \sum_{\text{minibatch}} \log \sum_{n \in \text{batch}} e^{-d(z, z_n)}$

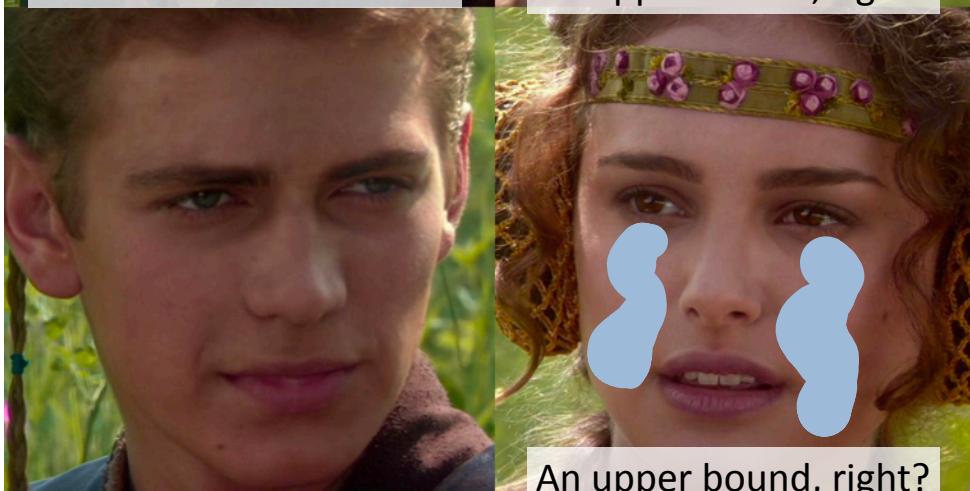
### Jensen's Inequality

$$\log \sum x \geq \sum \log x$$

$$\therefore d(\bar{z}, \bar{z}^+) + \log \sum_n e^{-d(z, z_n)} \geq d(\bar{z}, \bar{z}^+) + \sum_{\text{minibatch}} \log \sum_{n \in \text{batch}} e^{-d(z, z_n)}$$

: when we minimize over the minibatch objective, we're minimizing the LB (not the UB) on our original objective.  
 $\rightarrow$  Not good  $\rightarrow$  may not be minimizing on the original objective.

↳ Larger BS  $\rightarrow$  closer to minimizing on the original objective.



# Why does contrastive learning need a large batch size?

Interpretation of loss: classifying augmented example from rest of dataset

$$\mathcal{L}_{N\text{-way}}(\theta) = - \sum_i \log \frac{\exp(-d(\tilde{z}_i, \tilde{z}'_i))}{\sum_{j \neq i} \exp(-d(\tilde{z}_i, \tilde{z}_j))} \quad \leftarrow \text{summation over entire dataset}$$

Intuition: Closest  $z$  will dominate the denominator, can be missed when subsampling

Mathematically: Minimizing a lower-bound. 🎉

# Solutions to requiring a large batch size

## 1. Store representations from previous batches (“momentum contrast”)

He, Fan, Wu, Xie, Girshick. MoCo. CVPR 2020

- Good results with mini batch size of 256

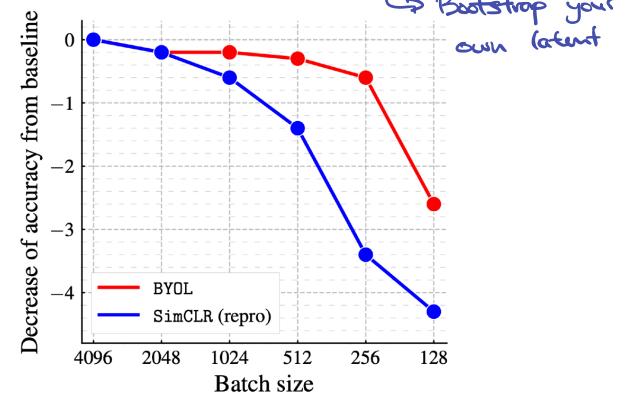
↳ Not exactly correct :: encoder will change over the course of training, ~ previous representations will be outdated.

↳ but still allows (some) decoupling of BS from the estimate

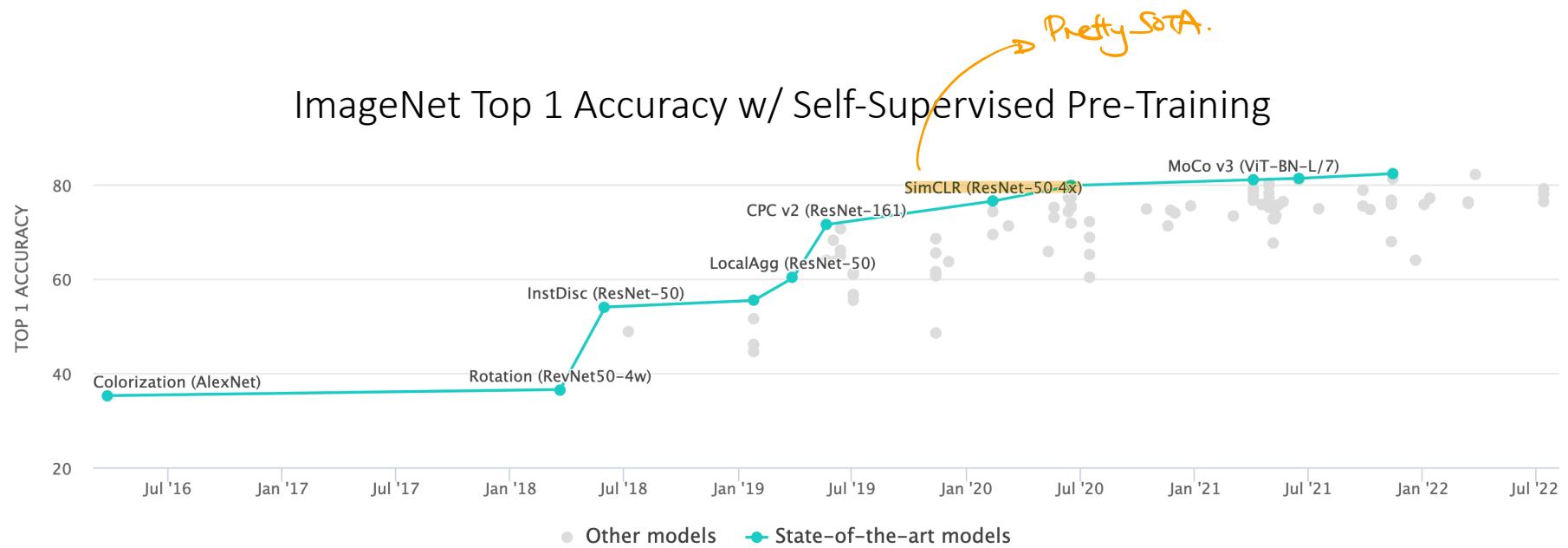
## 2. Predict representation of same image under different augmentation (“BYOL”)

Grill\*, Strub\*, Altché\*, Tallec\*, Richemond\*, et al. BYOL. NeurIPS 2020

- ❖ No negatives required!
- More resilient to batch size



# Performance of contrastive learning



Contrastive methods are near state-of-the-art in self-supervised pre-training for visual data.

# Contrastive learning beyond augmentations

We don't have good engineered augmentations for many applications!

## 1. *Learn* the augmentations in adversarial manner (but perturbations bounded to $\ell_1$ sphere)

Tamkin, Wu, Goodman. Viewmaker Networks. ICLR 2021

- competitive with SimCLR on image data
- good results on speech & sensor data

Interesting

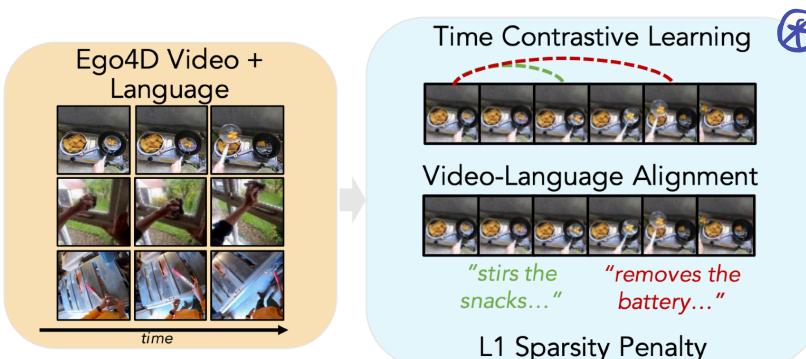
↳ Bounded to some amount  
(otherwise it'll give completely arbitrary images)



⊗ (+) yes : Frames closer in time  
⊗ yes : Frames further apart in time, or from other videos.

## 2. *Time-contrastive learning* on *videos* effective for robotics pre-training

Nair, Rajeswaran, Kumar, Finn, Gupta. R3M. CoRL 2022.



Given 20 demos (<10 min of supervision)



60% success



40% success

# Contrastive learning beyond augmentations

We don't have good engineered augmentations for many applications!

1. *Learn* the augmentations in adversarial manner (but perturbations bounded to  $\ell_1$  sphere)

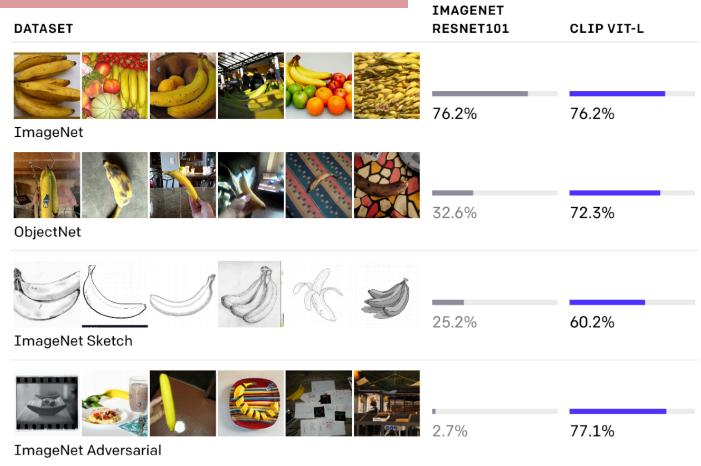
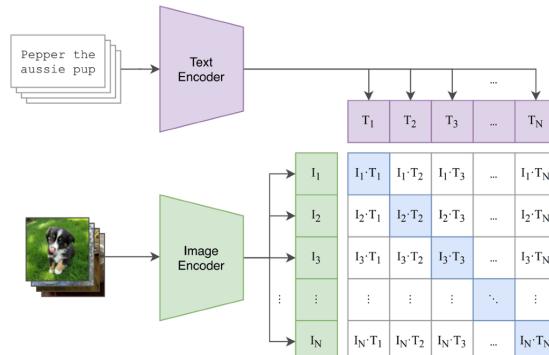
Tamkin, Wu, Goodman. Viewmaker Networks. ICLR 2021

2. *Time-contrastive learning* on *videos* effective for robotics pre-training

Nair, Rajeswaran, Kumar, Finn, Gupta. R3M. CoRL 2022.

3. **Image-text contrastive pre-training produces robust zero-shot models**

Radford\*, Kim\*, et al. CLIP. 2021.



# Summary of Contrastive Learning

## Pros:

- + General, effective framework
- + No generative modeling required  
↳ only need an encoder ( $f_\theta(x)$ )
- + Can incorporate domain knowledge through augmentations

 generate other (f)ues

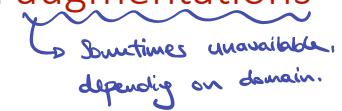
⊗ Means you can probs get away w/ a smaller model than if you had used a generative model.

## Challenges:

- *Negatives* can be hard to select
- Often requires *large batch size*
- Most successful with *augmentations*



 Smaller model, but need a large BS (generally)

 Sometimes unavailable, depending on domain.

# This Lecture

Unsupervised representation learning for few-shot learning

Part I: Contrastive learning

Part II (next time): Reconstruction-based methods

**Relation to meta-learning.**

# Contrastive Learning as Meta-Learning

## Meta-learning algorithm

1. Given unlabeled dataset  $\{x_i\}$ .
2. Create image class  $y_i$  from each datapoint via data augmentation  $\mathcal{D}_i := \{\tilde{x}_i, \tilde{x}'_i, \dots\}$
3. Run your favorite meta-learning algorithm.



### Differences:

- SimCLR samples **one task** per minibatch; meta-learning usually samples **multiple**
- SimCLR compares **all pairs** of samples; meta-learning compares query examples only to support examples & not to other query examples.

↳ SimCLR is more data-efficient w/ the batch

Otherwise the two learning methods are v. similar

# Contrastive Learning as Meta-Learning

## Meta-learning algorithm

1. Given unlabeled dataset  $\{x_i\}$ .
2. Create image class  $y_i$  from each datapoint via data augmentation  $\mathcal{D}_i := \{\tilde{x}_i, \tilde{x}'_i, \dots\}$
3. Run your favorite meta-learning algorithm.

Contrastive vs. meta-learning representations, transfer from ImageNet

	Flowers102	DTD	VOC2007	Aircraft	Food101	SUN397	CIFAR-10	CIFAR-100
Non-parametric meta learner ↗	SimCLR	92.4	72.7	66.0	83.7	86.3	57.4	94.8
	ProtoNet	92.7	71.5	64.7	83.9	86.2	56.4	96.0
	R2-D2	<b>94.5</b>	<b>73.8</b>	<b>69.9</b>	<b>86.2</b>	<b>86.9</b>	<b>59.7</b>	<b>96.7</b>

Optim-based  
meta learner ↗

Representations transfer similarly well.

[Few-Shot Learning]: Give it a few examples, embed, ~ meta clf. (@ test time)

[Contrastive learning]: Learn the representation, ~ fine-tune the whole network  
(@ test time)

(@ test time)

# Lecture Outline

Unsupervised representation learning for few-shot learning

Part I: Contrastive learning

Part II (next time): Reconstruction-based methods

^ next lecture by **TA Eric Mitchell**  
(NLP PhD student)

Relation to meta-learning.

## Goals for the lecture:

- Understand **contrastive learning**: intuition, design choices, how to implement
- How contrastive learning relates to meta-learning

# Course Reminders

Project proposal due Wednesday.  
(graded lightly, for your benefit)

Homework 2 due next Monday 10/24.